```matlab
 1 %% MACHINE LEARNING PROJECT
 2 % ALGORITHMS PERFORMING THE PREDICTION OF PRICES OF USED CARS (REGRESSION TASK)
 3
 4 % The two algorithms to be performed and compared are :
 5  % 1. Random Forrest Regressor
 6  % 2. Decision Tree for Regression
 7
 8  %% About the dataset
 9
10 % The data set contains 17,965 rows with 9 columns. The columns and their↙
descriptions are as follows:
11
12 %- model --> Model Name.
13 %- year --> The year it was bought.
14 %- price --> The price at which the used car will be sold.
15 %- transmission --> The transmission type i.e. manual, automatic or semi-automatic.
16 %- mileage --> The miles that used car has driven.
17 %- fuelType --> The fuel type of the car i.e. petrol, diesel, hybrid, electric or↙
other
18 %- tax --> The tax that will be applied on the selling price of that used car.
19 %- mpg --> The miles per gallon ratio telling us how many mies it can drive per↙
gallon of fuel.
20 %- engineSize --> The engine size of the used car.
21
22 %% Uploading the ford dataset from the 100,000 UK Used Car Data set from kaggle
23 clc
24 % reading my ford csv file & converting it into a table
25 fordcar = readtable("C:\Users\FIDVI FATEMA\OneDrive\Desktop\machine learning\10k↙
used car dataset\ford.csv");
26 %  displaying the head of the dataset
27 head(fordcar);
28
29 %% EXPLORATORY DATA ANALYSIS
30
31 %% 1. Checking Missing Values in my dataset
32 clc
33 missing_values = sum(ismissing(fordcar));
34
35 % Displaying the number of missing values for each variable
36 disp('Number of missing values for each variable:');
37 disp(missing_values);
38
39 %there are no missing values in my dataset
40
41 %% 2. Checking if there are any duplicates in my dataset
42 clc
43 [~, unique_indices, ~] = unique(fordcar);
44
45 % Getting the duplicate rows and removing them
46 duplicate_rows = fordcar(setdiff(1:height(fordcar), unique_indices), :);
47 fordcar_noduplicates = fordcar(unique_indices, :);
```

```matlab
48 % There were  154 duplicate rows
49
50 % Looking at the dimensions of the new fordcar_noduplicates dataset after the ↙
   removal of the duplilcates
51 dimensions = size(fordcar_noduplicates);
52 disp(['Number of rows: ', num2str(dimensions(1))]);
53 disp(['Number of columns: ', num2str(dimensions(2))]);
54
55 % After the removal of the duplicates, the shape of fordcar_final is (17811 x 9)
56
57 %% 3. Descriptive Summary of the variables in my dataset
58 clc
59 summary(fordcar_noduplicates)
60
61 % By doing this, we can see the minimum, maximum and the median values of
62 % each of the clumns in the dataset.
63
64 %% 4. Visually analyzing the distribution for price column
65
66 % Extracting the 'price' column
67 price = fordcar_noduplicates.price;
68 % Creating a histogram
69 figure;
70 histogram(price, 'BinWidth', 5000);
71 title('Distribution of Price');
72 xlabel('Price');
73 ylabel('Frequency');
74
75 % from this distribution graph, it is clear that  the overwhelming majority
76 % of price falls in the 0-25,000 range.
77
78 %% 5. Understanding the distribution of categorical columns
79 clc
80 % Extracting the 'fuelType' column
81 fuelType = fordcar_noduplicates.fuelType;
82 % Using the tabulate function to get fuelType counts
83 fuelType_counts = tabulate(fuelType);
84 % Displaying the results
85 disp(fuelType_counts);
86
87 % Performing the same for 'model' and 'transmission' columns
88 model = fordcar_noduplicates.model;
89 model_counts = tabulate(model);
90 disp(model_counts);
91
92 transmission = fordcar_noduplicates.transmission;
93 transmission_counts = tabulate(transmission);
94 disp(transmission_counts);
95
96 %% 6. Checking if there are outliers using boxplots
97
```

```matlab
 98 % creating a list for numerical columns
 99 num_list = {'price', 'mileage', 'tax', 'mpg'};
100
101 % Creating Boxplots for them
102 figure('Position', [100, 100, 900, 600]);
103 sgtitle('Distribution of Outliers', 'FontSize', 16);
104 rows = 2;
105 cols = 2;
106
107 % Iterating through numerical columns
108 for k = 1:numel(num_list)
109     % Creating subplot
110     subplot(rows, cols, k);
111
112     % Plot the boxplot
113     boxplot(fordcar_noduplicates.(num_list{k}));
114     title(num_list{k}, 'FontSize', 14);
115     if k > numel(num_list) - cols
116         xlabel('Variable');
117     end
118
119     if mod(k, cols) == 1
120         ylabel('Value');
121     end
122 end
123
124 %Even though there are a lot of outliers seen in the boxplots, removing them would⤶
not be the best option.
125 % For example, for price column, it is very much likely to have varied set of values⤶
ranging from a very high to very low values.
126 % Similarly, we can conclude for tax or mileage or mpg (mileage per gallon). If i⤶
remove them, it may break the flow of the dataset.
127 % Even though some of the data points can be actual outliers, I would conduct my⤶
research without removing them.
128
129 %% 7. Looking at the correlation between year and the sale price of the used cars
130
131 figure;
132 scatter(fordcar_noduplicates.year, fordcar_noduplicates.price, 'filled');
133 title('Ford cars sale price by year');
134 xlabel('Year');
135 ylabel('Price');
136
137 % From this, we can see that, the more recent the car is, the more pricier it is.
138 % There is one point that we can see as an outlier which shows the price of the car⤶
in 2060.
139 % This seems like an obvious outlier and should be removed from my dataset.
140
141 % removing the obvious outlier
142 logical_index = (fordcar_noduplicates.year ~= 2060);
143 fordcar_final = fordcar_noduplicates(logical_index, :);
```

```matlab
144
145 % Again looking at the distribution after removal of the outlier
146 figure;
147 scatter(fordcar_final.year, fordcar_final.price, 'filled');
148 title('Ford used car sale price by year');
149 xlabel('Year');
150 ylabel('Price');
151 %After removing this outlier, we can see the distribution of the price with respect
to years much more clearly.
152 % In this, we can see that, the more recent the car is i.e, the closer the year of
buying the car, the higher is the selling price of the car.
153 % That is, the lesser the age of the car, the higher is the price.
154
155 %% 8. Correlation matrix between my numeric variables
156
157 % Extracting the columns of interest
158 selected_columns_final = fordcar_final(:, {'price', 'year', 'mileage', 'tax', 'mpg',
'engineSize'});
159
160 % Converting the table to a matrix
161 selected_matrix = table2array(selected_columns_final);
162
163 % Calculating the correlation matrix
164 correlation_matrix_final = corrcoef(selected_matrix, 'rows', 'pairwise');
165
166 % Creating a new figure and heatmap
167 figure;
168 heatmap(selected_columns_final.Properties.VariableNames, selected_columns_final.
Properties.VariableNames, correlation_matrix_final, 'Colormap',bone, 'ColorLimits', [-1
1], 'FontSize', 10, 'CellLabelColor','k', 'CellLabelFormat', '%.2f');
169 title('Correlation Matrix (Pearson) for Numeric Values');
170 xlabel('Variables');
171 ylabel('Variables');
172
173 % This heatmap tells us a lot about our dataset and how our variables are correlated
to each other as well as with our target column i.e
174 % the price. The highest positive correlation with price column is 0.65 and it is
with 'year' column, while the strongest negative
175 % correlation with price column is -0.53 and it is with 'mileage' column. The
correlations of all numeric variables with my
176 % target column are quite high which is a good thing. The heatmap obtained from this
code only has the numerical columns and
177 % not the categorical columns of 'model', 'fuelType' and 'transmission'as they are
not label encoded yet.
178
179 %% 9. Distribution of the average price of the cars with the year column
180 clc
181 groupedData = groupsummary(fordcar_final, 'year', 'mean', 'price');
182 disp(groupedData);
183
184 %from this, we can see that the prices of the older cars i.e the cars from 1996-1998
```

```matlab
    are sold at a higher rates and then there is
185 %decline in the prices after those years as cars from after 1998 are not considered✔
    as antiques but just older and non-functional.
186 %This might justify why we see a positice correlation between 'price' and 'year'✔
    column.
187
188 %% 10. Distribution of the models of ford car in the dataset
189 clc
190 model_distribution = tabulate(fordcar_final.model);
191 disp(model_distribution);
192
193 %% 11. Correlation between model of the Ford car and the price
194 clc
195 % Extracting relevant columns
196 models = fordcar_final.model;
197 prices = fordcar_final.price;
198
199 % Creating a table with unique models and their corresponding total prices
200 unique_models = unique(models);
201 total_prices = zeros(size(unique_models));
202
203 for i = 1:length(unique_models)
204     total_prices(i) = sum(prices(strcmp(models, unique_models{i})));
205 end
206
207 % Creating a bar plot
208 figure;
209 bar(unique_models, total_prices);
210 title('Sum of Model Prices for Ford Cars');
211 xlabel('Ford Models');
212 ylabel('Total Price');
213 xtickangle(90);
214
215 % from this bar graph, it is very clear that there are 3 specific models of the ford✔
    company car namely 'Fiesta', 'Focus'
216 % & 'Kuga' which have a huge sum of prices i.e the sale price of these models might✔
    be higher than the others.
217 % the other reason might be that the numbers of these models in our datset is✔
    greater than the other models and so
218 % the total prices of these models exceed the others. The latter seems more accurate✔
    looking at the distribution of the model columns above.
219
220 %% Feature Engineering (Label Encoding)
221
222 %'Label Encoding is a technique that is used to convert categorical columns into✔
    numerical ones so that they can be fitted by machine
223 % learning models which only takes numerical data.
224 % for my dataset, we have 3 categorical columns namely 'transmission', 'fuelType'✔
    and 'model' and one 'year' column.
225
226 %% a) One-Hot Encoding :
```

```matlab
227 clc
228 % code help was taken from 'https://uk.mathworks.↙
com/help/deeplearning/ref/onehotencode.html'
229 % The 'transmission' column has 3 unique features while 'fuelType' column has 5↙
unique features which are not ordinal.
230 % So, the decision was made to apply one-hot encoding on these two categorical ↙
columns.
231
232
233 transmission = fordcar_final.transmission;
234 fuelType = fordcar_final.fuelType;
235
236 % Converting categorical variables to numerical indices
237 transmission_indices = grp2idx(categorical(fordcar_final.transmission));
238 fuelType_indices = grp2idx(categorical(fordcar_final.fuelType));
239
240 % Performing one-hot encoding
241 transmission_encoded = full(sparse(1:numel(transmission_indices),↙
transmission_indices, 1));
242 fuelType_encoded = full(sparse(1:numel(fuelType_indices), fuelType_indices, 1));
243
244 % Concatenating the encoded variables to the original table
245 fordcar_final_encoded = [fordcar_final, array2table(transmission_encoded,↙
'VariableNames', strcat('transmission_encoded_', cellstr(unique(fordcar_final.↙
transmission))))];
246 fordcar_final_encoded = [fordcar_final_encoded, array2table(fuelType_encoded,↙
'VariableNames', strcat('fuelType_encoded_', cellstr(unique(fordcar_final.fuelType))))];
247
248 % Displaying the head of the table
249 head(fordcar_final_encoded);
250
251 %% b) Integer Label Encoding
252
253 % This type of label encoding is the basic one as you can just assign an integer to↙
the specific categorical value in the column.
254 % The disadvantage of using this type of labelEncoder is that the model assumes that↙
the higher integer has more value than the
255 % lower one which is not true if the categorical variables in our column are not↙
ordinal.
256 % using this labelEncoding method for my 'model' column is more practical as the↙
model column has 23 unique categorical
257 % values and using one-hot encoding will just increase the dimensionality and make↙
my dataset more sparsed.
258
259 % Using integer label encoding on my 'model' column and assigning higher integer to↙
the more frequent models in my dataset.
260
261
262 % Getting the unique models and their counts
263 [uniqueModels, ~, modelIdx] = unique(fordcar_final_encoded.model);
264 % Getting the frequency of each model
```

```matlab
265 modelCounts = histcounts(modelIdx, 1:numel(uniqueModels)+1);
266 % Getting the rank of each model based on frequency
267 [~, rankedModels] = sort(modelCounts, 'descend');
268 % Assigning the ranked models as 'model_encoded'
269 fordcar_final_encoded.model_encoded = zeros(size(fordcar_final_encoded, 1), 1);
270
271 % Loop through uniqueModels and assigning corresponding rank
272 for i = 1:numel(uniqueModels)
273     modelIndices = ismember(fordcar_final_encoded.model, uniqueModels{i});
274     fordcar_final_encoded.model_encoded(modelIndices) = rankedModels(i);
275 end
276
277 % We can later see whether the results of the algorithms vary if we just
278 % assign the variables in 'model' columns with integers without taking into
279 % account the frequency of the model in our dataset.
280
281 %% c) Ordinal label encoding for my 'year' Column
282 % code help taken from 'https://uk.mathworks.com/help/matlab/matlab_prog/ordinal-
categorical-arrays.html'
283 % Performing Ordinal encoding for my year column as the recent the year is, the
higher the integer value it will be assigned.
284 fordcar_final_encoded.year = categorical(fordcar_final_encoded.year);
285 fordcar_final_encoded.year_encoded = grp2idx(fordcar_final_encoded.year);
286
287 % Dropping the original categorical columns
288 categorical_columns_to_remove = {'transmission', 'fuelType','model','year'};
289 fordcar_final_encoded = removevars(fordcar_final_encoded,
categorical_columns_to_remove);
290
291 % Displaying the head of the updated table
292 head(fordcar_final_encoded)
293
294 %  Dimensions of the dataset after encoding
295 dim = size(fordcar_final_encoded);
296 disp(['Number of rows: ', num2str(dim(1))]);
297 disp(['Number of columns: ', num2str(dim(2))]);
298
299 %%  Splitting the data into training, testing and validation (70:20:10 ratio)
300
301 % 70 is the training set ratio
302 % 20 is the testing set ratio
303 % 10 is the validation set ratio which we will use for performing Grid
304 % Search for Hyperparameter tuning for our models
305 clc
306 % Converting table to matrix
307 data = fordcar_final_encoded{:,:};
308
309 % Setting the random seed for reproducibility i.e each time we run this code, we get
the same results.
310 rng(43);
311
```

```matlab
312 % Create a partition for 70% training and 30% testing
313 cv = cvpartition(size(data, 1), 'HoldOut', 0.3);
314
315 % Training data
316 training_data = data(cv.training, :);
317
318 % Hold-out data (combined testing and validation i.e 20:10)
319 Holdoutdata = data(cv.test, :);
320
321 % Further splitting the hold-out data into testing (20%) and validation (10%)
322 cv2 = cvpartition(size(Holdoutdata, 1), 'HoldOut', 0.3333);
323
324 % Testing data
325 testing_data = Holdoutdata(cv2.training, :);
326
327 % Validation data
328 validation_data = Holdoutdata(cv2.test, :);
329
330 % Display the sizes of the datasets
331 disp(['Training Data Size: ' num2str(size(training_data, 1))]);      %70
332 disp(['Testing Data Size: ' num2str(size(testing_data, 1))]);        %20
333 disp(['Validation Data Size: ' num2str(size(validation_data, 1))]);  %10
334
335 %% Extracting the features and target variable for training, testing and validation
336
337 % as price is our second column
338 % Extracting features and target variable from validation_data
339 X_validation = validation_data(:, (2:end));
340 y_validation = validation_data(:, 1);
341
342 % Extracting features and target variable from training_data
343 X_train = training_data(:, (2:end));
344 y_train = training_data(:, 1);
345
346 % Extracting features and target variable from testing_data
347 X_test = testing_data(:, (2:end));
348 y_test = testing_data(:, 1);
349
350 %------------------------------------------------------------------------
351 %% RANDOM FORREST REGRESSION ALGORITHM
352
353 % Statistics and Machine Learning Toolbox installed
354
355 %% 1. Grid Search for finding the best Hyperparameters for Random Forest model
356
357 % Code help taken from 'https://uk.mathworks.com/help/stats/regression-tree-↙
ensembles.html'
358 clc
359 % setting random seed for reproducibility
360 rng(43);
361
```

```matlab
362 % Defining the range of hyperparameters to search over
363 numTreesRange = 10:50:300;
364 maxNumSplitsRange = 10:10:150;
365
366 % Initializing variables to store results
367 bestNumTrees = 0;
368 bestMaxNumSplits = 0;
369 rfbestR2 = -Inf;
370 rfbestMSE = Inf;
371
372 % Performing grid search
373 for numTrees = numTreesRange
374     for maxNumSplits = maxNumSplitsRange
375
376         % Creating Random Forest model for regression for grid search
377         rfmodel_gridsearch = TreeBagger(numTrees, X_train, y_train, 'Method', ↙
'regression', 'MaxNumSplits', maxNumSplits);
378         % Evaluating random forest on the validation set
379
380         y_predrf_validation = predict(rfmodel_gridsearch, X_validation);
381
382         % Calculating R-squared value
383         rfR2_validation = corr(y_predrf_validation, y_validation)^2;
384         % Calculating Mean Squared Error (MSE)
385         rfMSE_validation = mean((y_predrf_validation - y_validation).^2);
386
387     % Checking if the current hyperparameters result in better R-squared value and ↙
lower MSE
388         if rfR2_validation > rfbestR2 && rfMSE_validation < rfbestMSE
389             bestNumTrees = numTrees;
390             bestMaxNumSplits = maxNumSplits;
391             rfbestR2 = rfR2_validation;
392             rfbestMSE = rfMSE_validation;
393         end
394     end
395 end
396
397 % Displaying the best hyperparameters, R-squared and MSE value
398 disp(['Best Number of Trees for random forest: ' num2str(bestNumTrees)]);
399 disp(['Best MaxNumSplits for random forest: ' num2str(bestMaxNumSplits)]);
400 disp(['Best R-squared Value for random forest after gridsearch: ' num2str↙
(rfbestR2)]);
401 disp(['Best MSE Value for random forest after gridsearch: '  num2str(rfbestMSE)]);
402
403 % The computational time taken for this step is more than anticipated.
404 % The result of this grid search gives us the hyperparameters to get the best ↙
results for our prediction.
405
406 %% 2. Performing Random Forest Regression on my training data using the best ↙
hyperparameters from GridSearch
407 clc
```

```matlab
408 % using 60 no. of decision trees and 150 splits from grid search to get best
409 % results from training the model on training set.
410
411 % setting random seed for reproducibility
412 rng(43);
413
414 % Best hyperparameters from grid search
415 bestNumTrees = 60;
416 bestMaxNumSplits = 150;
417
418 % Training my random forest model on training set
419 randomforestmodel = TreeBagger(bestNumTrees, X_train, y_train, 'Method',↙
'regression', 'MaxNumSplits', bestMaxNumSplits, 'OOBPredictorImportance', 'on');
420
421 % save the best model for random forest
422 save("bestmodel for random forest.mat", 'randomforestmodel');
423
424 % Making predictions on the training set
425 y_predrf_train = predict(randomforestmodel, X_train);
426
427 % Calculating R-squared value and MSE value for my training model
428 rfMSE_train = mean((y_train - y_predrf_train).^2);
429 rfR2_train = 1 - sum((y_train - y_predrf_train).^2) / sum((y_train - mean(y_train)).↙
^2);
430 disp(['MSE Value for Random Forest Model : ' num2str(rfMSE_train)]);
431 disp(['R-squared Value for Random Forest Model: ' num2str(rfR2_train)]);
432
433 % Creating a table containing original and predicted prices
434 trainingprediction_rf = table(y_train, y_predrf_train);
435 trainingprediction_rf.Properties.VariableNames = {'Original Price', 'Predicted↙
Price'};
436
437 % Displaying the top 10 rows of the table
438 display(trainingprediction_rf(1:10, :));
439
440 % Using the bestNumTrees as 60 and bestMaxNumSplits as 150 from gridsearch, the↙
baseline model gives
441 % a high r-squared value of 0.9215 which means that 92.15% variance in our predicted↙
values are explained by the model.
442 % Although MSE value of 1751473.18 is huge but it depends on the scale of our data↙
which is varied.
443
444 %% 3. Visualizing the difference between Original price and the predicted price in↙
training data
445
446 % Creating a line plot for predictive price and actual price for training set
447 figure('Position', [100, 100, 800, 500]);
448 plot(1:length(y_train), y_train, 'r-', 'LineWidth', 0.2);
449 hold on;
450 plot(1:length(y_predrf_train), y_predrf_train, 'b-', 'LineWidth', 0.2);
451 grid on;
```

```matlab
452 xlabel('Data Point Index');
453 ylabel('Price');
454 legend('Actual Price', 'Predicted Price');
455 title('Actual vs. Predicted Prices for training set for Random Forest');
456
457 %% 4. Feature Importance selection for Random Forest
458 % code help taken from 'https://uk.mathworks.com/help/stats/select-predictors-for-↙
    random-forests.html'
459 rng(43);
460 % Extract feature importance scores
461 importance_scores = randomforestmodel.OOBPermutedVarDeltaError;
462 % Specify the number of top features to select
463 num_top_features = 5;
464
465 % Find indices of top N features
466 [~, top_feature_indices] = maxk(importance_scores, num_top_features);
467 % Filter datasets to include only top features
468 X_train_selected = X_train(:, top_feature_indices);
469 X_validation_selected = X_validation(:, top_feature_indices);
470 X_test_selected = X_test(:, top_feature_indices);
471
472 % Retrain Random Forest on the selected features
473 randomforestmodel_selected = TreeBagger(bestNumTrees, X_train_selected, y_train,↙
    'Method', 'regression', 'MaxNumSplits', bestMaxNumSplits);
474
475 % Make predictions from the selected columns on the training set
476 y_predrf_train_selected = predict(randomforestmodel_selected, X_train_selected);
477
478 % Calculate R-squared value and MSE value for my training model
479 rfMSE_train_selected = mean((y_train - y_predrf_train_selected).^2);
480 rfR2_train_selected = 1 - sum((y_train - y_predrf_train_selected).^2) / sum((y_train↙
    - mean(y_train)).^2);
481
482 % Display R-squared value
483 disp(['MSE Value for Random Forest Model for selected columns : ' num2str↙
    (rfMSE_train_selected)]);
484 disp(['R-squared Value for Random Forest Model for selected models: ' num2str↙
    (rfR2_train_selected)]);
485
486 % While Feature selection does seem like an important step for predictive
487 % modeling, its application to my model yielded minimal impact on performance.
488 % Notably, the R-squared (r2) value remained consistent, and though there
489 % was a modest reduction in Mean Squared Error (MSE), it was not substantial
490 % enough to alter the overall model outcomes. This suggests that the initially
491 % chosen features already provided relevant information for prediction,
492 % and the removal of some features did not yield a significant improvement.
493 % I prefer having a model with all features for interpretability, even if it
494 % comes at the cost of a slightly more complex model.
495
496 %% 5. Evaluating the baseline model on the test data (unseen data)
497 clc
```

```matlab
498 % evaluating the performance of random forest model on the unseen test data without ↙
the feature selection.
499 % setting random seed for reproducibility
500 rng(43);
501
502 load("bestmodel for random forest.mat");
503
504 % Evaluating the baseline model on the testing set
505 y_predrf_test = predict(randomforestmodel, X_test);
506
507 % Calculating MSE and R-squared values on the testing set
508 rfMSE_test = mean((y_test - y_predrf_test).^2);
509 rfR2_test = 1 - sum((y_test - y_predrf_test).^2) / sum((y_test - mean(y_test)).^2);
510 disp(['MSE on Testing Set for Random Forest Model: ' num2str(rfMSE_test)]);
511 disp(['R-squared Value on Testing Set for Random Forest Model: ' num2str↙
(rfR2_test)]);
512
513 % R2 value for my test set or unseen data is 0.91799 which is considered outstanding ↙
but the MSE value
514 % of this model on my test set is 1828169.73. The lower the MSE value, the better ↙
the result is, but our
515 % MSE value depends on the variability of my data and I want to check whether this ↙
high mse value is for
516 % the entire data or just the test set via k-fold cross validation.
517
518 %% 6. Visualizing the difference between Original price and the predicted price in ↙
testing set
519
520 % Creating a line plot for predictive price vs. actual price on test set
521 figure('Position', [100, 100, 800, 500]);
522 plot(1:length(y_test), y_test, 'g-', 'LineWidth', 0.2);
523 hold on;
524 plot(1:length(y_predrf_test), y_predrf_test, 'm-', 'LineWidth', 0.2);
525 grid on;
526 xlabel('Data Point Index');
527 ylabel('Price');
528 legend('Actual Price', 'Predicted Price');
529 title('Actual vs. Predicted Prices for test set for Random Forest');
530
531 %% 7. Visualizing the predictive and actual prices for the training and testing data
532
533 % Creating a line plot for collectively seeing the performance of the model on ↙
training as well as testing set
534 figure('Position', [100, 200, 1100, 500]);
535 plot(y_train, 'r', 'LineWidth', 0.2, 'DisplayName', 'Actual Training');
536 hold on;
537 plot(y_predrf_train, 'b', 'LineWidth', 0.2, 'DisplayName', 'Fitted Training');
538 plot(length(y_train) + (1:length(y_test)), y_test, 'g', 'LineWidth', 0.2,↙
'DisplayName', 'Actual Testing');
539 plot(length(y_train) + (1:length(y_test)), y_predrf_test, 'm', 'LineWidth', 0.2,↙
'DisplayName', 'Fitted Testing');
```

```matlab
540 grid on;
541
542 title('RANDOM FOREST MODEL PREDICTIONS (on training and test set)');
543 xlabel('Observation');
544 ylabel('Price');
545 legend('show');
546
547 %% 8. Performing K-fold Cross Validation on my entire dataset for Random Forest
548 clc
549 % Performing K-fold Cross Validation after running my model on the training
550 % as well as testing data is crucial for robust model evaluation.
551 % It provides a more reliable estimate of the model's performance,
552 % reducing the risk of overfitting or underfitting observed with a single train-test↙
split.
553
554 % setting random seed for reproducibility
555 rng(43);
556 % Combine the training, validation, and testing sets
557 combined_data = [training_data; validation_data; testing_data];
558
559 %extract the features to put in X and Y
560 X = combined_data(:, (2:end));
561 y = combined_data(:,1);
562
563 % Define k-fold partition
564 cvpart = cvpartition(size(X, 1), 'kfold', 10);
565
566 % Performing k-fold cross-validation
567 mse_scores_rf = [];
568 r2_scores_rf = [];
569 for i = 1:cvpart.NumTestSets
570     % Extract training and testing data
571     trainIdx = cvpart.training(i);
572     testIdx = cvpart.test(i);
573     X_train = X(trainIdx, :);
574     X_test = X(testIdx, :);
575     y_train = y(trainIdx);
576     y_test = y(testIdx);
577
578     % Getting my random forest model
579     randomforestmodel = TreeBagger(bestNumTrees, X_train, y_train, 'Method',↙
'regression', 'MaxNumSplits', bestMaxNumSplits);
580
581     % Evaluating the model
582     y_predrf = predict(randomforestmodel, X_test);
583
584     % Calculating MSE
585     kfMSE_rf = mean((y_predrf - y_test).^2);
586     mse_scores_rf = [mse_scores_rf; kfMSE_rf];
587
588     % Calculating R-squared value
```

```matlab
589     kfR2_rf = corr(y_predrf, y_test)^2;
590     r2_scores_rf = [r2_scores_rf; kfR2_rf];
591
592     % Displaying results for each fold
593     disp(['Fold ' num2str(i) ' - MSE: ' num2str(kfMSE_rf) ', R-squared: ' num2str↙
(kfR2_rf)]);
594 end
595
596 % Analyzing overall results
597 average_mse_rf = mean(mse_scores_rf);
598 average_r2_rf = mean(r2_scores_rf);
599 disp(['Average MSE for k-fold: ' num2str(average_mse_rf)]);
600 disp(['Average R-squared for k-fold: ' num2str(average_r2_rf)]);
601
602 %The results obtained from k-fold cross-validation for our Random Forest model show↙
promising performance.
603 % The model effectively captures a substantial portion of the variability in the↙
target variable.
604 % The consistent MSE values indicate that overfitting is not a concern; instead, the↙
observed higher
605 % MSE values are suggestive of the inherent complexity or dispersion within the↙
dataset.
606 % It is important to recognize that these elevated MSE values do not necessarily↙
reflect poor model
607 % performance but rather highlight the nuanced characteristics of the data,↙
signaling potential complexity or variability.
608
609 %------------------------------------------------------------------------
610 %%  DECISION TREE REGRESSION ALGORITHM
611
612 % For decision tree, I will be using the split data and the extracted features and↙
target variables
613 % for training, testing and validation that I used earlier for random forest↙
regression model.
614
615 %% 1. Grid Search to find the best Hyperparameters for Decision Tree Model
616
617 % code help taken from 'https://uk.mathworks.com/help/stats/decision-trees.↙
html#bsw6p3v'
618 % Performing Grid Search method for getting the best Hyperparameters for Decision↙
Tree Model.
619 clc
620 % Set the random seed for reproducibility
621 rng(43);
622
623 % Defining the range of hyperparameters to search over
624 maxsplitsrange = 5:5:100;
625 minLeafSizeRange = 5:5:50;
626
627 % Initializing variables to store results
628 bestMaxSplits = 0;
```

```matlab
629 bestMinLeafSize = 0;
630 dtbestR2 = -Inf;
631 dtbestMSE = Inf;
632
633 % Performing grid search on validation set
634 for maxnumsplits = maxsplitsrange
635     for minLeafSize = minLeafSizeRange
636
637         % Creating decision tree model for regression for gridsearch
638         dtreemodel_gridsearch = fitrtree(X_train, y_train, 'MinLeafSize',↙
minLeafSize, 'MaxNumSplits', maxnumsplits);
639
640         % Evaluating the decision tree on the validation set
641         y_preddt_validation = predict(dtreemodel_gridsearch, X_validation);
642
643         % Calculating R-squared value
644         dtR2_validation = corr(y_preddt_validation, y_validation)^2;
645         % Calculating Mean Squared Error (MSE)
646         dtMSE_validation = mean((y_preddt_validation - y_validation).^2);
647
648         % Checking if the current hyperparameters result in better R-squared value↙
and lower MSE
649         if dtR2_validation > dtbestR2 && dtMSE_validation < dtbestMSE
650             bestMaxSplits = maxnumsplits;
651             bestMinLeafSize = minLeafSize;
652             dtbestR2 = dtR2_validation;
653             dtbestMSE = dtMSE_validation;
654         end
655     end
656 end
657
658 % Displaying the best hyperparameters and performance metrics
659 disp(['Best MaxNumber of splits: ' num2str(bestMaxSplits)]);
660 disp(['Best MinLeafSize: ' num2str(bestMinLeafSize)]);
661 disp(['Best R-squared Value for Decision Tree Model: ' num2str(dtbestR2)]);
662 disp(['Best MSE Value for Decision Tree Model: ' num2str(dtbestMSE)]);
663
664 % Computational time taken for gridsearch for decision tree is way lesser than that↙
taken for random forest.
665 % The result of this hyperparameter tuning gives us the best number of splits and↙
the minimum leaf size so as
666 % to obtain the best predictive model by using these hyperparameters.
667
668 %% 2. Performing  Decision Tree Model on my training set using the best↙
hyperparameters from GridSearch
669 clc
670 % to get the best results for our decision tree's model on our dataset, we'll be↙
using the hyperparameters
671 % we got from the grid search for our Decision Trees model
672
673 % setting random seed for reproducibility
```

```matlab
674 rng(43);
675
676 % Best hyperparameters from grid search
677 bestMaxSplits = 100;
678 bestMinLeafSize = 20;
679
680 % Training decision tree model for regression on training set
681 dtreemodel = fitrtree(X_train, y_train, 'MinLeafSize', bestMinLeafSize, ↵
'MaxNumSplits', bestMaxSplits);
682
683 %save the best model for decision tree
684 save("bestmodel for decision tree.mat", "dtreemodel");
685
686 % Make predictions on the training set
687 y_preddt_train = predict(dtreemodel, X_train);
688
689 % Calculate MSE and R2 values for the decision tree model
690 dtMSE_train = mean((y_train - y_preddt_train).^2);
691 dtR2_train = 1 - sum((y_train - y_preddt_train).^2) / sum((y_train - mean(y_train)). ↵
^2);
692 disp(['R-squared Value for Decision Tree Model on training set: ' num2str↵
(dtR2_train)]);
693 disp(['MSE Value for Decision Tree Model on training set: ' num2str(dtMSE_train)]);
694
695 % Creating a table containing original and predicted prices for the decision tree ↵
model
696 trainingprediction_dt = table(y_train, y_preddt_train);
697 trainingprediction_dt.Properties.VariableNames = {'Original Price', 'Predicted ↵
Price'};
698
699 % Displaying the top 10 rows of the table for the decision tree model
700 display(trainingprediction_dt(1:10, :));
701
702 %% 3. Visualizing the difference between Original price and the predicted price in↵
training data for Decision Tree
703
704 % Creating a line plot for predictive price and actual price for training set for↵
the decision tree model
705 figure('Position', [100, 100, 800, 500]);
706 plot(1:length(y_train), y_train, 'r-', 'LineWidth', 0.2);
707 hold on;
708 plot(1:length(y_preddt_train), y_preddt_train, 'b-', 'LineWidth', 0.2);
709 grid on;
710 xlabel('Data Point Index');
711 ylabel('Price');
712 legend('Actual Price', 'Predicted Price');
713 title('Actual vs. Predicted Prices for training set for Decision Tree');
714
715 %% 4. Evaluating the baseline model on the test data (unseen data) for Decision Tree
716 clc
717 % setting random seed for reproducibility
```

```matlab
718 rng(43);
719
720 load("bestmodel for decision tree.mat");
721
722 % Evaluating the decision tree model on the testing set
723 y_preddt_test = predict(dtreemodel, X_test);
724
725 % Calculating MSE and R-squared values on the testing set
726 dtMSE_test = mean((y_test - y_preddt_test).^2);
727 dtR2_test = 1 - sum((y_test - y_preddt_test).^2) / sum((y_test - mean(y_test)).^2);
728
729 % Displaying the MSE and R-squared values on the testing set for the decision tree ↙
model
730 disp(['MSE on Testing Set (Decision Tree Model): ' num2str(dtMSE_test)]);
731 disp(['R-squared Value on Testing Set (Decision Tree Model): ' num2str(dtR2_test)]);
732
733 % It is performing as good on the test set as it did on the training set.
734
735 %% 5. Visualizing the difference between Original price and the predicted price in ↙
testing set for Decision Tree
736
737 % Creating a line plot for predictive price vs. actual price on test set for ↙
decision tree
738 figure('Position', [100, 100, 800, 500]);
739 plot(1:length(y_test), y_test, 'g-', 'LineWidth', 0.2);
740 hold on;
741 plot(1:length(y_preddt_test), y_preddt_test, 'm-', 'LineWidth', 0.2);
742 grid on;
743 xlabel('Data Point Index');
744 ylabel('Price');
745 legend('Actual Price', 'Predicted Price (Decision Tree)');
746 title('Actual vs predicted price for test set for Decision Tree Model');
747
748 %% 6. Visualizing the predictive and actual prices for the training and testing data ↙
for decision tree model
749
750 % Creating a line plot for collectively seeing the performance of the model on ↙
training as well as testing set
751 figure('Position', [100, 200, 1100, 500]);
752 plot(y_train, 'r', 'LineWidth', 0.2, 'DisplayName', 'Actual Training');
753 hold on;
754 plot(y_preddt_train, 'b', 'LineWidth', 0.2, 'DisplayName', 'Fitted Training');
755 plot(length(y_train) + (1:length(y_test)), y_test, 'g', 'LineWidth', 0.2, ↙
'DisplayName', 'Actual Testing');
756 plot(length(y_train) + (1:length(y_test)), y_preddt_test, 'm', 'LineWidth', 0.2, ↙
'DisplayName', 'Fitted Testing');
757 grid on;
758
759 title('DECISION TREE MODEL PREDICTIONS (on training and testing set)');
760 xlabel('Observation');
761 ylabel('Price');
```

```matlab
762 legend('show');
763
764    %% 7. Performing K-fold Cross Validation on my entire dataset
765
766 % K-fold cross validation is done to make sure that my model is not overfitting the↙
training data.
767 % It has the same motive of providing a more reliable estimate of the model's↙
generalization performance,
768 % helping identify potential overfitting or underfitting issues and ensuring the↙
model's consistency across diverse data subset.
769 clc
770 % setting random seed for reproducibility
771 rng(43);
772
773 % Using the combined data and X and y features that we created above when performing↙
the K-fold cross validation on random forest model.
774
775 % Defining k-fold partition
776 cvpart = cvpartition(size(X, 1), 'kfold', 10);
777
778 % Perform k-fold cross-validation on decision tree model (following the same↙
procedure as for random forest)
779 mse_scores_dt = [];
780 r2_scores_dt = [];
781 for i = 1:cvpart.NumTestSets
782     % Extracting training and testing data
783     trainIdx = cvpart.training(i);
784     testIdx = cvpart.test(i);
785     X_train = X(trainIdx, :);
786     X_test = X(testIdx, :);
787     y_train = y(trainIdx);
788     y_test = y(testIdx);
789
790   % Getting my decision tree model
791   dtreemodel = fitrtree(X_train, y_train, 'MinLeafSize', bestMinLeafSize,↙
'MaxNumSplits', bestMaxSplits);
792
793     % Evaluating the model
794     y_preddt = predict(dtreemodel, X_test);
795
796     % Calculating MSE
797     kfMSE_dt = mean((y_preddt - y_test).^2);
798     mse_scores_dt = [mse_scores_dt; kfMSE_dt];
799
800     % Calculating R-squared value
801     kfR2_dt = corr(y_preddt, y_test)^2;
802     r2_scores_dt = [r2_scores_dt; kfR2_dt];
803
804     % Displaying results for each fold
805     disp(['Fold ' num2str(i) ' - MSE: ' num2str(kfMSE_dt) ', R-squared: ' num2str↙
(kfR2_dt)]);
```

```matlab
806 end
807
808 % Analyzing overall results
809 average_mse_dt = mean(mse_scores_dt);
810 average_r2_dt = mean(r2_scores_dt);
811 disp(['Average MSE for k-fold for Decision Tree model: ' num2str(average_mse_dt)]);
812 disp(['Average R-squared for k-fold for Decision Tree model: ' num2str↙
(average_r2_dt)]);
813
814 %--------------------------------------------------------------------------
815 %% Q. Does scaling the data before applying Random Forest and Decision Tree models↙
yield comparable results to unscaled data, or does it significantly impact model↙
performance?
816 %(extra analytical question)
817
818 % - Scaling of my variables (Through minmax scaling)
819
820 % Before scaling
821 figure;
822 subplot(2, 1, 1);
823 histogram(X_train(:, 1), 'BinEdges', linspace(min(X_train(:, 1)), max(X_train(:,↙
1)), 20));
824 title('Before Scaling - Feature 1');
825
826 % Calculate min and max values for scaling
827 min_values_train = min(X_train, [], 1);
828 max_values_train = max(X_train, [], 1);
829
830 % Scale training data
831 for feature_index = 1:size(X_train, 2)
832     X_train(:, feature_index) = (X_train(:, feature_index) - min_values_train↙
(feature_index)) / (max_values_train(feature_index) - min_values_train(feature_index));
833 end
834
835 % Scale validation data using the same min and max values
836 for feature_index = 1:size(X_validation, 2)
837     X_validation(:, feature_index) = (X_validation(:, feature_index) -↙
min_values_train(feature_index)) / (max_values_train(feature_index) - min_values_train↙
(feature_index));
838 end
839
840 % Scale test data using the same min and max values
841 for feature_index = 1:size(X_test, 2)
842     X_test(:, feature_index) = (X_test(:, feature_index) - min_values_train↙
(feature_index)) / (max_values_train(feature_index) - min_values_train(feature_index));
843 end
844
845 % After scaling
846 subplot(2, 1, 2);
847 histogram(X_train(:, 1), 'BinEdges', linspace(0, 1, 20));
848 title('After Scaling - Feature 1');
```

```matlab
849
850 % The subplots are just to check whether my scaling worked or not.
851
852 %%  Random Forest with scaled data
853 % 1.Grid Search
854 % setting random seed for reproducibility
855 rng(43);
856
857 % Defining the range of hyperparameters to search over
858 numTreesRange = 10:50:300;
859 maxNumSplitsRange = 10:10:150;
860
861 % Initializing variables to store results
862 bestNumTrees = 0;
863 bestMaxNumSplits = 0;
864 rfbestR2 = -Inf;
865 rfbestMSE = Inf;
866
867 % Performing grid search
868 for numTrees = numTreesRange
869     for maxNumSplits = maxNumSplitsRange
870
871         % Creating Random Forest model for regression
872         randomforestmodel = TreeBagger(numTrees, X_train, y_train, 'Method', ↙
'regression', 'MaxNumSplits', maxNumSplits);
873
874         % Evaluating random forest on the validation set
875         y_predrf = predict(randomforestmodel, X_validation);
876
877         % Calculating R-squared value
878         rfR2 = corr(y_predrf, y_validation)^2;
879          % Calculating Mean Squared Error (MSE)
880         rfMSE = mean((y_predrf - y_validation).^2);
881
882  % Checking if the current hyperparameters result in better R-squared value and↙
lower MSE
883         if rfR2 > rfbestR2 && rfMSE < rfbestMSE
884             bestNumTrees = numTrees;
885             bestMaxNumSplits = maxNumSplits;
886             rfbestR2 = rfR2;
887             rfbestMSE = rfMSE;
888         end
889     end
890 end
891
892 % Display the best hyperparameters and R-squared value
893 disp(['Best Number of Trees: ' num2str(bestNumTrees)]);
894 disp(['Best MaxNumSplits: ' num2str(bestMaxNumSplits)]);
895 disp(['Best R-squared Value for random forest: ' num2str(rfbestR2)]);
896 disp(['Best MSE Value for random forest:  '  num2str(rfbestMSE)]);
897
```

```matlab
898 % 2. Running the model on training set for Random Forest (Scaled Data)
899
900 % setting random seed for reproducibility
901 rng(43);
902
903 % Best hyperparameters from grid search
904 bestNumTrees = 60;
905 bestMaxNumSplits = 140;
906
907 % Training my random forest model on training set
908 randomforestmodel = TreeBagger(bestNumTrees, X_train, y_train, 'Method',↙
'regression', 'MaxNumSplits', bestMaxNumSplits);
909
910 % Make predictions on the training set
911 y_predrf_train = predict(randomforestmodel, X_train);
912
913 % Calculate R-squared value and MSE value for my training model
914 rfMSE_train = mean((y_train - y_predrf_train).^2);
915 rfR2_train = 1 - sum((y_train - y_predrf_train).^2) / sum((y_train - mean(y_train)).↙
^2);
916
917 % Display R-squared value
918 disp(['MSE Value for Random Forest Model : ' num2str(rfMSE_train)]);
919 disp(['R-squared Value for Random Forest Model: ' num2str(rfR2_train)]);
920
921 % There isn't much of a difference between these results and the ones
922 % obtained without scaling
923
924 %% Decision Tree with scaled data
925 % 1. Grid Search
926 % Set the random seed for reproducibility
927 rng(43)
928
929 % Define the range of hyperparameters to search over
930 maxsplitsrange = 5:5:250;
931 minLeafSizeRange = 5:5:50;
932
933 % Initialize variables to store results
934 bestMaxSplits = 0;
935 bestMinLeafSize = 0;
936 dtbestR2 = -Inf;
937 dtbestMSE = Inf;
938
939 % Performing grid search on validation set
940 for maxnumsplits = maxsplitsrange
941     for minLeafSize = minLeafSizeRange
942
943         % Creating decision tree model for regression
944         dtreemodel = fitrtree(X_train, y_train, 'MinLeafSize', minLeafSize,↙
'MaxNumSplits', maxnumsplits);
945
```

```matlab
946          % Evaluating the decision tree on the validation set
947          y_preddt = predict(dtreemodel, X_validation);
948
949          % Calculating R-squared value
950          dtR2 = corr(y_preddt, y_validation)^2;
951
952          % Calculating Mean Squared Error (MSE)
953          dtMSE = mean((y_preddt - y_validation).^2);
954
955          % Checking if the current hyperparameters result in better R-squared value↵
and lower MSE
956          if dtR2 > dtbestR2 && dtMSE < dtbestMSE
957              bestMaxSplits = maxnumsplits;
958              bestMinLeafSize = minLeafSize;
959              dtbestR2 = dtR2;
960              dtbestMSE = dtMSE;
961          end
962      end
963 end
964
965 % Display the best hyperparameters and performance metrics
966 disp(['Best MaxNumber of splits: ' num2str(bestMaxSplits)]);
967 disp(['Best MinLeafSize: ' num2str(bestMinLeafSize)]);
968 disp(['Best R-squared Value for Decision Tree Model: ' num2str(dtbestR2)]);
969 disp(['Best MSE Value for Decision Tree Model: ' num2str(dtbestMSE)]);
970
971 % 2. Running the model on training set for Decision Tree (Scaled Data)
972
973 % setting random seed for reproducibility
974 rng(43);
975
976 % Best hyperparameters from grid search
977 bestMaxSplits = 100;
978 bestMinLeafSize = 20;
979
980 % Creating decision tree model for regression
981 dtreemodel = fitrtree(X_train, y_train, 'MinLeafSize', bestMinLeafSize,↵
'MaxNumSplits', bestMaxSplits);
982
983 % Make predictions on the training set
984 y_preddt_train = predict(dtreemodel, X_train);
985
986 % Calculate MSE value for the decision tree model
987 dtMSE_train = mean((y_train - y_preddt_train).^2);
988
989 % Calculate R-squared value for the decision tree model
990 dtR2_train = 1 - sum((y_train - y_preddt_train).^2) / sum((y_train - mean(y_train)).↵
^2);
991
992 % Display the calculated values
993 disp(['R-squared Value for Decision Tree Model: ' num2str(dtR2_train)]);
```

```matlab
 994 disp(['MSE Value for Decision Tree Model: ' num2str(dtMSE_train)]);
 995
 996 % I am getting the same values as I got before scaling for grid research for↙
decision tree model.
 997
 998 % Scaling the data had minimal impact on the results of Random Forest, while having↙
no impact on
 999 % Decision Tree Model's performance . Consequently, it can be inferred that both↙
Random Forest and Decision Tree algorithms
1000 % exhibit robustness to non-scaled data, performing admirably without the need for↙
normalization or scaling.
1001 %-------------------------------------------------------------------------
```