

▼ Importing Data

```
1 from google.colab import drive
2 drive.mount('/content/drive')
```

🔗 Drive already mounted at /content/drive; to attempt to forcibly remount, call drive.mount("/content/drive", force_remount=True).

```
1 #IMPORTING THE DATA FROM THE DATASETS.
2 import pandas as pd
3 import numpy as np
4 import csv
5 import matplotlib.pyplot as plt
6 from skimage import data, io, filters
7 # Import the necessary modules and libraries
8 import numpy as np
9 from sklearn.tree import DecisionTreeRegressor
10 import matplotlib.pyplot as plt
11 from sklearn.model_selection import KFold
12
13 path3 = "/content/drive/My Drive/Homework 4/OnlineNewsPopularityTest.csv"
14 path = "/content/drive/My Drive/Homework 4/OnlineNewsPopularityTrain.csv"
15
16 train_data_df = pd.read_csv(path)
17 train_data_df = train_data_df.drop('url', axis=1)
18 train_data = train_data_df.drop(' shares', axis=1)
19 feature_list = np.array(train_data.columns)
20 train_outcomes= train_data_df[' shares']
21
22 test_data_df = pd.read_csv(path3)
23 test_data_df = test_data_df.drop('url', axis=1)
24 test_data = test_data_df.drop(' shares', axis=1)
25 test_outcomes= test_data_df[' shares']
```

-

```

1 X = np.array(train_data)
2 y = np.array(train_outcomes)
3 kf = KFold()
4 kf.get_n_splits(X)

```

↳ 5

```

1 def mse (y1, y2):
2     summation = 0 #variable to store the summation of differences
3     n = len(y1) #finding total number of items in list
4     for i in range (n): #looping through each element of the list
5         difference = y1[i] - y2[i] #finding the difference between observed and predicted value
6         squared_difference = difference**2 #taking square of the difference
7         summation = summation + squared_difference #taking a sum of all the differences
8     MSE = summation/n
9     return(MSE)
10

```

▼ Decision Tree Regression

```

1 # Fit regression model
2 k_folds = np.array_split(train_data_df,5)
3 tree_depth = []
4 model_score = []
5 for j in range (1,10):
6     mean_square_error= []
7     for k in range(5):
8         train = k_folds.copy()
9         validation = k_folds [k]
10        validation_y = validation[' shares']
11        validation_y = np.array(validation_y)
12        validation = validation.drop(' shares', axis =1)
13        validation = np.array(validation)
14        del train[k]
15        train = pd.concat(train, sort =False)

```

```

16 train_y = train[' shares']
17 train_y = np.array(train_y)
18 train = train.drop (' shares',axis=1)
19 train = np.array(train)
20 regr_1 = DecisionTreeRegressor(criterion="mse",max_depth=j,random_state=0)
21 regr_1.fit(train, train_y)
22 y_1 = regr_1.predict(validation)
23 mean_square_error.append(mse(y_1, validation_y))
24 model_score.append(((sum(mean_square_error))/(len(mean_square_error))))
25 i=(np.argmin(model_score))
26 print("Average Scores for different Depths of tree: ", model_score)
27 print("Best Score: {0} at a tree depth of: {1}".format(model_score[i],[i]))

```

➞ Average Scores for different Depths of tree: [1016.3140978011204, 894.0392691925359, 1050.7181523765987, 943.7684687466854, 766.0052357596807856 at a tree depth of: [5]

```

1 print(np.argmin(model_score))

```

➞ 5

```

1 regr_2 = DecisionTreeRegressor(criterion="mse",max_depth=i,random_state=0)
2 regr_2.fit(X,y)
3 y_2 = regr_2.predict(test_data)
4 mean_square_error=(mse(y_2,test_outcomes))
5 print(mean_square_error)

```

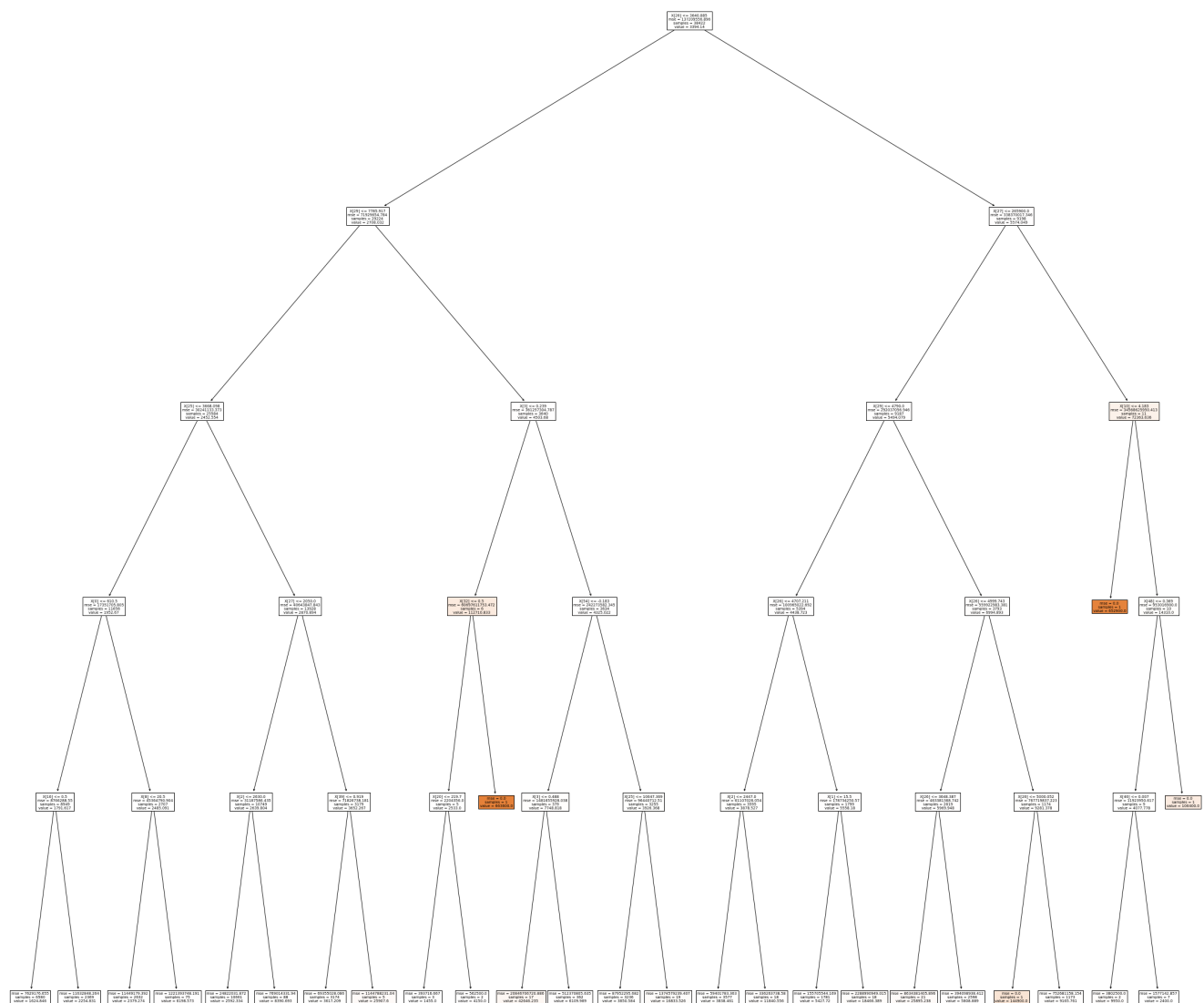
➞ 2611.1073794518447

```

1 #PLOTING THIS JUST FOR FUN.
2 from sklearn import tree
3 plt.figure(figsize=(40,40))
4 #tree.plot_tree(regr_2.fit(test_data, test_outcomes))
5 tree.plot_tree(regr_2, filled=True)
6 plt.show()

```

➞



▼ Random Forest Regression

```
1 # Fit regression model
2 def tree_regressor ():
3     data_in = train_data_df.sample(frac=0.3,axis=0)
4     k_folds = np.array_split(data_in,5)
5     y = data_in[' shares']
6     y = np.array(y)
7     X = data_in.drop(' shares', axis=1)
8     X = np.array(X)
9     tree_depth = []
10    model_score = []
11    y_output = []
12    for j in range (1,11):
13        mean_square_error= []
14        predicted_output = []
15        for k in range(5):
```

```

16     train = k_folds.copy()
17     validation = k_folds[k]
18     validation_y = validation['shares']
19     validation_y = np.array(validation_y)
20     validation = validation.drop('shares', axis=1)
21     validation = np.array(validation)
22     del train[k]
23     train = pd.concat(train, sort=False)
24     train_y = train['shares']
25     train_y = np.array(train_y)
26     train = train.drop('shares', axis=1)
27     train = np.array(train)
28     regr_1 = DecisionTreeRegressor(criterion="mse", max_depth=j, max_features="log2")
29     regr_1.fit(train, train_y)
30     y_1 = regr_1.predict(validation)
31     mean_square_error.append(mse(y_1, validation_y))
32     model_score.append(((sum(mean_square_error))/(len(mean_square_error))))
33     i=(np.argmin(model_score))
34     regr_out = DecisionTreeRegressor(criterion="mse", max_depth=j, max_features="log2")
35     regr_out.fit(X, y)
36     y_out = regr_out.predict(test_data)
37     return(y_out)

```

```

1 def random_forest (no_of_trees):
2     y_outcomes = []
3     for k in range (1, no_of_trees):
4         y_outcomes.append(tree_regressor())
5     summed_outcomes = np.zeros(len(y_outcomes[0]))
6     for i in range(len(y_outcomes)):
7         for j in range(len(y_outcomes[i])):
8             summed_outcomes[j] = summed_outcomes[j] + y_outcomes[i][j]
9     summed_outcomes = summed_outcomes/(len(y_outcomes))
10    result = (mse(summed_outcomes, test_outcomes))
11    return(result)

```

```

1 best_forest = []
2 for i in range (2 10):

```

```
2 for i in range (2,10):
3     temp = (random_forest(i))
4     best_forest.append(temp)
5 i=np.argmin(best_forest)
6 print("Best performer error is {0} using a forest of {1} trees".format(best_forest[i],i+1))
```

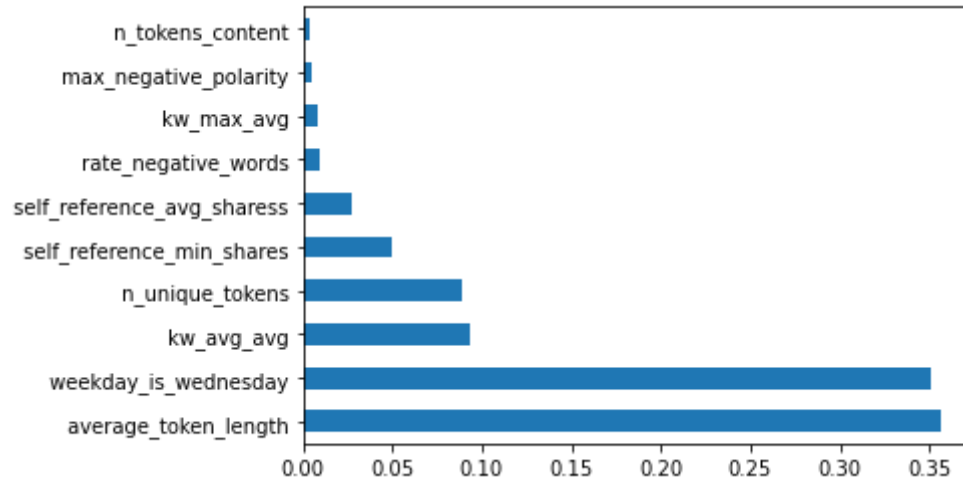
➞ Best performer error is 1477.6864650270134 using a forest of 8 trees

▼ Feature Exploration

```
1 print(regr_2.feature_importances_)
2 #plot graph of feature importances for better visualization
3 feat_importances = pd.Series(regr_2.feature_importances_,index=feature_list)
4 feat_importances.nlargest(10).plot(kind='barh')
5 plt.show()
```

➞


```
[9.62466339e-04 2.91732094e-03 3.93012428e-03 8.86356188e-02
0.00000000e+00 0.00000000e+00 0.00000000e+00 0.00000000e+00
1.02432754e-03 0.00000000e+00 3.56996536e-01 0.00000000e+00
0.00000000e+00 0.00000000e+00 0.00000000e+00 0.00000000e+00
6.65716402e-04 0.00000000e+00 0.00000000e+00 0.00000000e+00
8.39288447e-06 0.00000000e+00 0.00000000e+00 0.00000000e+00
0.00000000e+00 8.21802410e-03 9.29751009e-02 4.97306617e-02
0.00000000e+00 2.69224268e-02 0.00000000e+00 0.00000000e+00
3.50689021e-01 0.00000000e+00 0.00000000e+00 0.00000000e+00
0.00000000e+00 0.00000000e+00 0.00000000e+00 2.38854951e-03
0.00000000e+00 0.00000000e+00 0.00000000e+00 0.00000000e+00
0.00000000e+00 0.00000000e+00 8.53870264e-05 0.00000000e+00
9.07391896e-03 0.00000000e+00 0.00000000e+00 0.00000000e+00
0.00000000e+00 0.00000000e+00 4.77640717e-03 0.00000000e+00
0.00000000e+00 0.00000000e+00 0.00000000e+00]
```



Could not connect to the reCAPTCHA service. Please check your internet connection and reload to get a reCAPTCHA challenge.