```
1   import numpy as np
2   import pandas as pd
3   import matplotlib.cm as cm
4   import matplotlib.pyplot as plt
5   # Importing data
6   path = "" # Put path of your folder of your data if it's not in the same folder
7   data_train = pd.read_csv("train.csv")
8   data_train.head()
9   data_train.shape
```
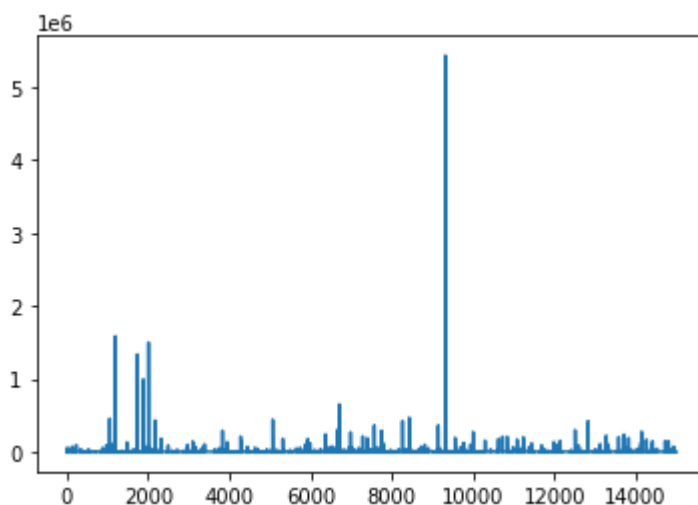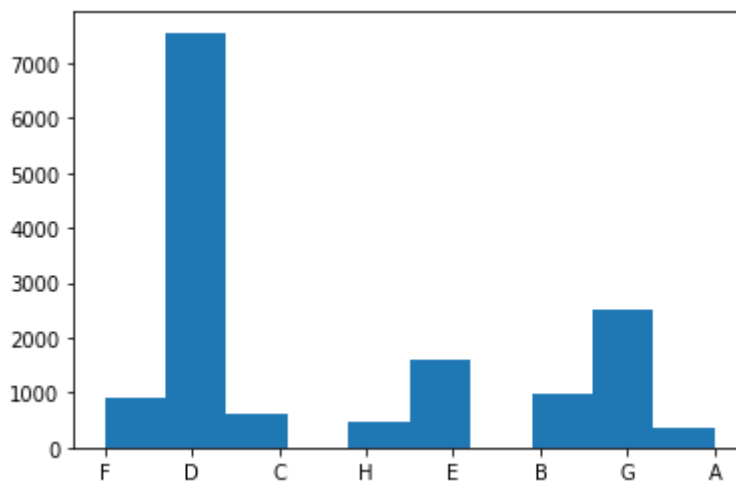
```
(14999, 9)
```

```
1   # Visualization
2   # Individual Plots
3   plt.hist(data_train["category"])
4   plt.show()
5   plt.plot(data_train["adview"])
6   plt.show()
7   # Remove videos with adview greater than 2000000 as outlier
8   data_train = data_train[data_train["adview"] <2000000]
```
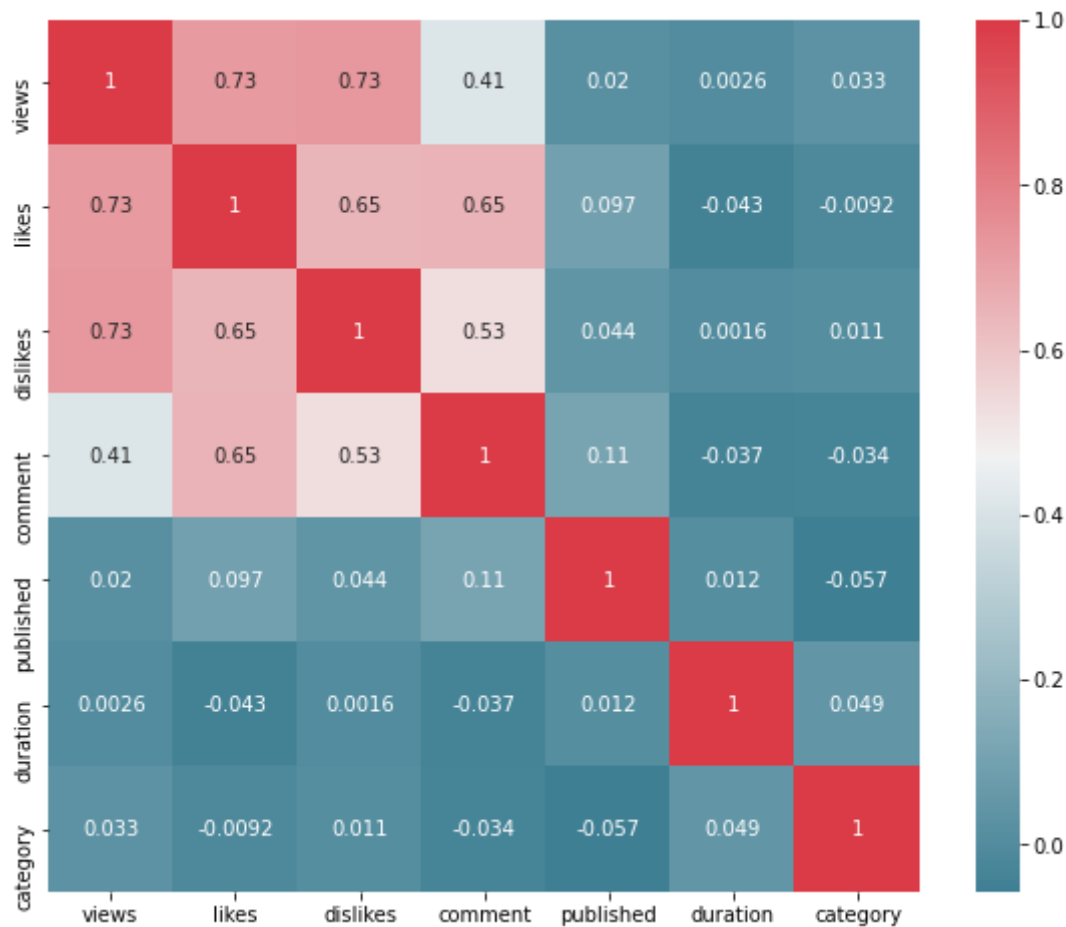


```
1   # Heatmap
2   import seaborn as sns
3   f, ax = plt.subplots(figsize=(10, 8))
4   corr = data_train.corr()
```

```
5  sns.heatmap(corr, mask=np.zeros_like(corr, dtype=np.bool), cmap=sns.diverging_palette
6  square=True, ax=ax,annot=True)
7  plt.show()
```



```
1   # Removing character "F" present in data
2   data_train=data_train[data_train.views!='F']
3   data_train=data_train[data_train.likes!='F']
4   data_train=data_train[data_train.dislikes!='F']
5   data_train=data_train[data_train.comment!='F']
6   data_train.head()
7   # Assigning each category a number for Category feature
8   category={'A': 1,'B':2,'C':3,'D':4,'E':5,'F':6,'G':7,'H':8}
9   data_train["category"]=data_train["category"].map(category)
10  data_train.head()
11
```

|   | vidid | adview | views | likes | dislikes | comment | published | duration | categor |
|---|-------|--------|-------|-------|----------|---------|-----------|----------|---------|
| 0 | VID_18655 | 40 | 1031602 | 8523 | 363 | 1095 | 2016-09-14 | PT7M37S | |
| 1 | VID_14135 | 2 | 1707 | 56 | 2 | 6 | 2016-10-01 | PT9M30S | |
| 2 | VID_2187 | 1 | 2023 | 25 | 0 | 2 | 2016-07-02 | PT2M16S | |
|   |  |  |  |  |  |  | 2016-07- | | |

```
1   # Convert values to integers for views, likes, comments, dislikes and adview
```

```
2   data_train["views"] = pd.to_numeric(data_train["views"])
3   data_train["comment"] = pd.to_numeric(data_train["comment"])
4   data_train["likes"] = pd.to_numeric(data_train["likes"])
5   data_train["dislikes"] = pd.to_numeric(data_train["dislikes"])
6   data_train["adview"]=pd.to_numeric(data_train["adview"])
7   column_vidid=data_train['vidid']
```

```
1   # Endoding features like Category, Duration, Vidid
2   from sklearn.preprocessing import LabelEncoder
3   data_train['duration']=LabelEncoder().fit_transform(data_train['duration'])
4   data_train['vidid']=LabelEncoder().fit_transform(data_train['vidid'])
5   data_train['published']=LabelEncoder().fit_transform(data_train['published'])
6   data_train.head()
```

|   | vidid | adview | views | likes | dislikes | comment | published | duration | category |
|---|-------|--------|-------|-------|----------|---------|-----------|----------|----------|
| 0 | 5912 | 40 | 1031602 | 8523 | 363 | 1095 | 2168 | 2925 | 6 |
| 1 | 2741 | 2 | 1707 | 56 | 2 | 6 | 2185 | 3040 | 4 |
| 2 | 8138 | 1 | 2023 | 25 | 0 | 2 | 2094 | 1863 | 3 |
| 3 | 9004 | 6 | 620860 | 777 | 161 | 153 | 2119 | 2546 | 8 |
| 4 | 122 | 1 | 666 | 1 | 0 | 0 | 2091 | 1963 | 4 |

```
1   # Convert Time_in_sec for duration
2   import datetime
3   import time
4   def checki(x):
5     y = x[2:]
6     h = ''
7     m = ''
8     s = ''
9     mm = ''
10    P = ['H','M','S']
11    for i in y:
12      if i not in P:
13        mm+=i
14      else:
15          if(i=="H"):
16           h = mm
17           mm = ''
18          elif(i == "M"):
19           m = mm
20           mm = ''
21          else:
22            s = mm
23            mm = ''
24    if(h==''):
25        h = '00'
26    if(m == ''):
27        m = '00'
28    if(s==''):
29        s='00'
```

```
30      bp = h+':'+m+':'+s
31      return bp
32  train=pd.read_csv("train.csv")
33  mp = pd.read_csv(path + "train.csv")["duration"]
34  time = mp.apply(checki)
35
36  def func_sec(time_string):
37      h, m, s = time_string.split(':')
38      return int(h) * 3600 + int(m) * 60 + int(s)
39
40  time1=time.apply(func_sec)
41
42  data_train["duration"]=time1
43  data_train.head()
```

|   | vidid | adview | views | likes | dislikes | comment | published | duration | category |
|---|-------|--------|-------|-------|----------|---------|-----------|----------|----------|
| 0 | 5912  | 40     | 1031602 | 8523 | 363      | 1095    | 2168      | 457      | 6        |
| 1 | 2741  | 2      | 1707  | 56    | 2        | 6       | 2185      | 570      | 4        |
| 2 | 8138  | 1      | 2023  | 25    | 0        | 2       | 2094      | 136      | 3        |
| 3 | 9004  | 6      | 620860 | 777  | 161      | 153     | 2119      | 262      | 8        |
| 4 | 122   | 1      | 666   | 1     | 0        | 0       | 2091      | 31       | 4        |

```
1  # Split Data
2  Y_train = pd.DataFrame(data = data_train.iloc[:, 1].values, columns = ['target'])
3  data_train=data_train.drop(["adview"],axis=1)
4  data_train=data_train.drop(["vidid"],axis=1)
5  data_train.head()
```

|   | views | likes | dislikes | comment | published | duration | category |
|---|-------|-------|----------|---------|-----------|----------|----------|
| 0 | 1031602 | 8523 | 363      | 1095    | 2168      | 457      | 6        |
| 1 | 1707  | 56    | 2        | 6       | 2185      | 570      | 4        |
| 2 | 2023  | 25    | 0        | 2       | 2094      | 136      | 3        |
| 3 | 620860 | 777  | 161      | 153     | 2119      | 262      | 8        |
| 4 | 666   | 1     | 0        | 0       | 2091      | 31       | 4        |

```
1  from sklearn.model_selection import train_test_split
2  X_train, X_test, y_train, y_test = train_test_split(data_train, Y_train, test_size=0.
3  X_train.shape
4  # Normalise Data
5  from sklearn.preprocessing import MinMaxScaler
6  scaler = MinMaxScaler()
7  X_train=scaler.fit_transform(X_train)
8  X_test=scaler.fit_transform(X_test)
9  X_train.mean()
```

```
0.1739096800320488
```

```
1    # Evaluation Metrics
2    from sklearn import metrics
3    def print_error(X_test, y_test, model_name):
4        prediction = model_name.predict(X_test)
5        print('Mean Absolute Error:', metrics.mean_absolute_error(y_test, prediction))
6        print('Mean Squared Error:', metrics.mean_squared_error(y_test, prediction))
7        print('Root Mean Squared Error:', np.sqrt(metrics.mean_squared_error(y_test, pred
```

```
1    # Linear Regression
2    from sklearn import linear_model
3    linear_regression = linear_model.LinearRegression()
4    linear_regression.fit(X_train, y_train)
5    print_error(X_test,y_test, linear_regression)
6
```

```
Mean Absolute Error: 3707.378005824532
Mean Squared Error: 835663131.1210337
Root Mean Squared Error: 28907.83857573986
```

```
1    # Support Vector Regressor
2    from sklearn.svm import SVR
3    supportvector_regressor = SVR()
4    supportvector_regressor.fit(X_train,y_train)
5    print_error(X_test,y_test, linear_regression)
6
```

```
/usr/local/lib/python3.7/dist-packages/sklearn/utils/validation.py:760: DataConversic
  y = column_or_1d(y, warn=True)
Mean Absolute Error: 3707.378005824532
Mean Squared Error: 835663131.1210337
Root Mean Squared Error: 28907.83857573986
```

```
1    # Decision Tree Regressor
2    from sklearn.tree import DecisionTreeRegressor
3    decision_tree = DecisionTreeRegressor()
4    decision_tree.fit(X_train, y_train)
5    print_error(X_test,y_test, decision_tree)
6
```

```
Mean Absolute Error: 2683.9286202185795
Mean Squared Error: 897836055.1950136
Root Mean Squared Error: 29963.91254818058
```

```
1    # Random Forest Regressor
2    from sklearn.ensemble import RandomForestRegressor
3    n_estimators = 400
4    max_depth = 50
5    min_samples_split= 30
6    min_samples_leaf= 3
7    random_forest = RandomForestRegressor(n_estimators = n_estimators, max_depth = max_de
8    random_forest.fit(X_train,y_train)
9    print_error(X_test,y_test, random_forest)
```

```
9    print_error(X_test,y_test, random_forest)
```

```
/usr/local/lib/python3.7/dist-packages/ipykernel_launcher.py:8: DataConversionWarning

Mean Absolute Error: 3416.1310229756986
Mean Squared Error: 742092301.0656787
Root Mean Squared Error: 27241.371130427313
```

```python
1    # Artificial Neural Network
2    import keras
3    from keras.layers import Dense
4    ann = keras.models.Sequential([
5                                    Dense(6, activation="relu",
6                                    input_shape=X_train.shape[1:]),
7                                    Dense(6,activation="relu"),
8                                    Dense(1)
9                                    ])
10   optimizer=keras.optimizers.Adam()
11   loss=keras.losses.mean_squared_error
12   ann.compile(optimizer=optimizer,loss=loss,metrics=["mean_squared_error"])
13   history=ann.fit(X_train,y_train,epochs=100)
14   ann.summary()
15   print_error(X_test,y_test,ann)
16
```

```
Epoch 1/100
366/366 [==============================] - 17s 3ms/step - loss: 469571112.9589 - m
Epoch 2/100
366/366 [==============================] - 1s 3ms/step - loss: 703724705.7231 - me
Epoch 3/100
366/366 [==============================] - 1s 3ms/step - loss: 1227204391.2371 - m
Epoch 4/100
366/366 [==============================] - 1s 3ms/step - loss: 1118232067.2933 - m
Epoch 5/100
366/366 [==============================] - 1s 3ms/step - loss: 661635794.7248 - me
Epoch 6/100
366/366 [==============================] - 1s 3ms/step - loss: 532626680.9373 - me
Epoch 7/100
366/366 [==============================] - 1s 3ms/step - loss: 952438887.0191 - me
Epoch 8/100
366/366 [==============================] - 1s 3ms/step - loss: 1011212140.5443 - m
Epoch 9/100
366/366 [==============================] - 1s 3ms/step - loss: 823545630.3433 - me
Epoch 10/100
366/366 [==============================] - 1s 3ms/step - loss: 794369160.1761 - me
Epoch 11/100
366/366 [==============================] - 1s 3ms/step - loss: 447305384.2603 - me
Epoch 12/100
366/366 [==============================] - 1s 3ms/step - loss: 660571687.4934 - me
Epoch 13/100
366/366 [====================-=========] - 1s 3ms/step - loss: 476336505.6621 - me
Epoch 14/100
366/366 [==========-===================] - 1s 3ms/step - loss: 761525071.2071 - me
Epoch 15/100
366/366 [=============.================] - 1s 3ms/step - loss: 524057261.6291 - me
Epoch 16/100
366/366 [==============================] - 1s 3ms/step - loss: 1686341799.4114 - m
```

```
Epoch 17/100
366/366 [==============================] - 1s 3ms/step - loss: 608129630.1580 - me
Epoch 18/100
366/366 [==============================] - 1s 3ms/step - loss: 363924695.4537 - me
Epoch 19/100
366/366 [==============================] - 1s 3ms/step - loss: 1045351705.5082 - m
Epoch 20/100
366/366 [==============================] - 1s 3ms/step - loss: 428724058.7420 - me
Epoch 21/100
366/366 [==============================] - 1s 3ms/step - loss: 442341584.7129 - me
Epoch 22/100
366/366 [==============================] - 1s 3ms/step - loss: 407823120.8604 - me
Epoch 23/100
366/366 [==============================] - 1s 3ms/step - loss: 1071861235.8072 - m
Epoch 24/100
366/366 [==============================] - 1s 3ms/step - loss: 581302252.1308 - me
Epoch 25/100
366/366 [==============================] - 1s 3ms/step - loss: 503853939.9108 - me
Epoch 26/100
366/366 [==============================] - 1s 3ms/step - loss: 998389941.9510 - me
Epoch 27/100
366/366 [==============================] - 1s 3ms/step - loss: 743672269.1172 - me
Epoch 28/100
366/366 [==============================] - 1s 3ms/step - loss: 538569881.5228 - me
Epoch 29/100
366/366 [                              ]   1s 3ms/step   loss: 012507401 8806   m
```

```python
1   #Saving Scikitlearn models
2   import joblib
3   joblib.dump(decision_tree, "decisiontree_youtubeadviewpred.pkl")
4   # Saving Keras Artificial Neural Network model
5   ann.save("ann_youtubeadviewpred.h5")
```

```python
1   import numpy as np
2   import pandas as pd
3   import matplotlib.cm as cm
4   import matplotlib.pyplot as plt
5   # Importing data
6   data_test = pd.read_csv("test.csv")
7   data_test.head()
8   data_test.shape
```

```
(8764, 8)
```

```python
1    # Removing character "F" present in data
2    data_test=data_test[data_test.views!='F']
3    data_test=data_test[data_test.likes!='F']
4    data_test=data_test[data_test.dislikes!='F']
5    data_test=data_test[data_test.comment!='F']
6    data_test.head()
7    # Assigning each category a number for Category feature
8    category={'A': 1, 'B':2, 'C':3, 'D':4, 'E':5, 'F':6, 'G':7, 'H':8}
9    data_test["category"]=data_test["category"].map(category)
10   data_test.head()
```

|   | vidid | views | likes | dislikes | comment | published | duration | category |
|---|-------|-------|-------|----------|---------|-----------|----------|----------|
| 0 | VID_1054 | 440238 | 6153 | 218 | 1377 | 2017-02-18 | PT7M29S | 2 |
| 1 | VID_18629 | 1040132 | 8171 | 340 | 1047 | 2016-06-28 | PT6M29S | 6 |
| 2 | VID_13967 | 28534 | 31 | 11 | 1 | 2014-03-10 | PT37M54S | 4 |
| 3 | VID_19442 | 1316715 | 2284 | 250 | 274 | 2010-06-05 | PT9M55S | 7 |
| 4 | VID_770 | 1893173 | 2519 | 225 | 116 | 2016-09-03 | PT3M8S | 2 |

```python
# Convert values to integers for views, likes, comments, dislikes and adview
data_test["views"] = pd.to_numeric(data_test["views"])
data_test["comment"] = pd.to_numeric(data_test["comment"])
data_test["likes"] = pd.to_numeric(data_test["likes"])
data_test["dislikes"] = pd.to_numeric(data_test["dislikes"])
column_vidid=data_test['vidid']
```

```python
# Endoding features like Category, Duration, Vidid
from sklearn.preprocessing import LabelEncoder
data_test['duration']=LabelEncoder().fit_transform(data_test['duration'])
data_test['vidid']=LabelEncoder().fit_transform(data_test['vidid'])
data_test['published']=LabelEncoder().fit_transform(data_test['published'])
data_test.head()
```

|   | vidid | views | likes | dislikes | comment | published | duration | category |
|---|-------|-------|-------|----------|---------|-----------|----------|----------|
| 0 | 231 | 440238 | 6153 | 218 | 1377 | 2053 | 2115 | 2 |
| 1 | 3444 | 1040132 | 8171 | 340 | 1047 | 1825 | 2055 | 6 |
| 2 | 1593 | 28534 | 31 | 11 | 1 | 1009 | 1506 | 4 |
| 3 | 3775 | 1316715 | 2284 | 250 | 274 | 116 | 2265 | 7 |
| 4 | 7644 | 1893173 | 2519 | 225 | 116 | 1892 | 1625 | 2 |

```python
# Convert Time_in_sec for duration
import datetime
import time
def checki(x):
    y = x[2:]
    h = ''
    m = ''
    s = ''
    mm = ''
    P = ['H','M','S']
    for i in y:
      if i not in P:
         mm+=i
      else:
          if(i=="H"):
            h = mm
            mm = ''
          elif(i == "M"):
            m = mm
```

```
19                    m = mm
20                    mm = ''
21                else:
22                    s = mm
23                    mm = ''
24        if(h==''):
25            h = '00'
26        if(m == ''):
27            m = '00'
28        if(s==''):
29            s='00'
30        bp = h+':'+m+':'+s
31        return bp
32    train=pd.read_csv("test.csv")
33    mp = pd.read_csv("test.csv")["duration"]
34    time = mp.apply(checki)
35
36    def func_sec(time_string):
37        h, m, s = time_string.split(':')
38        return int(h) * 3600 + int(m) * 60 + int(s)
39
40    time1=time.apply(func_sec)
41
42    data_test["duration"]=time1
43    data_test.head()
```

|   | vidid | views | likes | dislikes | comment | published | duration | category |
|---|-------|-------|-------|----------|---------|-----------|----------|----------|
| 0 | 231 | 440238 | 6153 | 218 | 1377 | 2053 | 449 | 2 |
| 1 | 3444 | 1040132 | 8171 | 340 | 1047 | 1825 | 389 | 6 |
| 2 | 1593 | 28534 | 31 | 11 | 1 | 1009 | 2274 | 4 |
| 3 | 3775 | 1316715 | 2284 | 250 | 274 | 116 | 595 | 7 |
| 4 | 7644 | 1893173 | 2519 | 225 | 116 | 1892 | 188 | 2 |

```
1    # Split Data
2    Z_train = pd.DataFrame(data = data_test.iloc[:, 1].values, columns = ['target'])
3    data_test=data_test.drop(["vidid"],axis=1)
4    data_test.head()
```

|   | views | likes | dislikes | comment | published | duration | category |
|---|-------|-------|----------|---------|-----------|----------|----------|
| 0 | 440238 | 6153 | 218 | 1377 | 2053 | 449 | 2 |
| 1 | 1040132 | 8171 | 340 | 1047 | 1825 | 389 | 6 |
| 2 | 28534 | 31 | 11 | 1 | 1009 | 2274 | 4 |
| 3 | 1316715 | 2284 | 250 | 274 | 116 | 595 | 7 |
| 4 | 1893173 | 2519 | 225 | 116 | 1892 | 188 | 2 |

```
1    from sklearn.model_selection import train_test_split
2    X_train, X_test, y_train, y_test = train_test_split(data_test, Z_train, test_size=0.2
```

```
3   X_train.shape
4   # Normalise Data
5   from sklearn.preprocessing import MinMaxScaler
6   scaler = MinMaxScaler()
7   X_train=scaler.fit_transform(X_train)
8   X_test=scaler.fit_transform(X_test)
9   X_train.mean()
```

```
0.1691684916653104
```

```
1   # Evaluation Metrics
2   from sklearn import metrics
3   def print_error(X_test, y_test, model_name):
4       prediction = model_name.predict(X_test)
5       print('Mean Absolute Error:', metrics.mean_absolute_error(y_test, prediction))
6       print('Mean Squared Error:', metrics.mean_squared_error(y_test, prediction))
7       print('Root Mean Squared Error:', np.sqrt(metrics.mean_squared_error(y_test, pred
8
```

```
1   # Decision Tree Regressor
2   from sklearn.tree import DecisionTreeRegressor
3   decision_tree = DecisionTreeRegressor()
4   decision_tree.fit(X_train, y_train)
5   print_error(X_test,y_test, decision_tree)
```

```
Mean Absolute Error: 3458894.505263158
Mean Squared Error: 115794316199804.33
Root Mean Squared Error: 10760776.747047786
```

```
1   #Saving Scikitlearn models
2   import joblib
3   joblib.dump(decision_tree, "decisiontree_prediction.pkl")
```

```
['decisiontree_prediction.pkl']
```

```
1   from sklearn.tree import DecisionTreeRegressor
2   regressor = DecisionTreeRegressor()
3   regressor.fit(X_train, y_train)
```

```
DecisionTreeRegressor(ccp_alpha=0.0, criterion='mse', max_depth=None,
                      max_features=None, max_leaf_nodes=None,
                      min_impurity_decrease=0.0, min_impurity_split=None,
                      min_samples_leaf=1, min_samples_split=2,
                      min_weight_fraction_leaf=0.0, presort='deprecated',
                      random_state=None, splitter='best')
```

```
1   y_pred = regressor.predict(X_test)
```

```
1   df=pd.DataFrame({'Predictions':y_pred})
2   df.head()
```

|   | Predictions |
|---|---|
| 0 | 265021.0 |
| 1 | 14179.0 |
| 2 | 9682759.0 |
| 3 | 24248747.0 |
| 4 | 425196.0 |

```
1    predictions=df
```

```
1    predictions
```

|      | Predictions |
|------|-------------|
| 0    | 265021.0 |
| 1    | 14179.0 |
| 2    | 9682759.0 |
| 3    | 24248747.0 |
| 4    | 425196.0 |
| ...  | ... |
| 1705 | 13572428.0 |
| 1706 | 199992.0 |
| 1707 | 6011924.0 |
| 1708 | 5662900.0 |
| 1709 | 1005358.0 |

1710 rows × 1 columns

```
1    predictions.to_csv('predictions.csv', index=False)
```

```
1    ! cat predictions.csv
```

```
Predictions
265021.0
14179.0
9682759.0
24248747.0
425196.0
356072.0
350410.0
4283780.0
800112.0
28556617.0
3046514.0
6947.0
14957150.0
```

```
27112851.0
29390.0
878075.0
45510.0
543526.0
891984.0
203763.0
430245.0
4006843.0
7841531.0
276.0
259759.0
4401353.0
5767561.0
49727.0
5646790.0
44952.0
1385406.0
3725457.0
292963.0
6376460.0
1333017.0
1621910.0
87814.0
1409323.0
2983.0
1368706.0
356503.0
483380.0
11483.0
1568927.0
216077.0
567505.0
16539356.0
1143033.0
185502.0
1504436.0
327481.0
542673.0
740999.0
1339.0
1511094.0
468635.0
166598.0
1718530.0
```

1