

```
1 import numpy as np
2 import pandas as pd
3 import matplotlib.pyplot as plt
4 %matplotlib inline
5
6 import datetime
7
8 from sklearn.model_selection import train_test_split
9
10 from sklearn.tree import DecisionTreeClassifier
11 from sklearn.metrics import accuracy_score
12 from sklearn.metrics import precision_score
13 from sklearn.metrics import jaccard_score
14 from sklearn.metrics import balanced_accuracy_score
15
16 from sklearn.preprocessing import StandardScaler
17
18 from sklearn.tree import DecisionTreeRegressor
19 from sklearn.ensemble import RandomForestRegressor
20
21 from sklearn.metrics import mean_squared_error
22 from sklearn.metrics import mean_absolute_error
23 from math import sqrt
24
```

```
1 url_country = 'https://raw.githubusercontent.com/Advait177013/CSE4020_ML_JComp/main/clean
2 url_income = 'https://raw.githubusercontent.com/Advait177013/CSE4020_ML_JComp/main/clean
3 url_region = 'https://raw.githubusercontent.com/Advait177013/CSE4020_ML_JComp/main/clean
4 df_country = pd.read_csv(url_country)
5 df_income = pd.read_csv(url_income)
6 df_region = pd.read_csv(url_region)
```

Team Members :

Advait Deochakke

Harish T R

Shah Siddh

```
1 #Code by Advait
2
3 #creating a copy to facilitate operations
4 df=df_country.copy()
5
6 #changing the date format to a format which we can easily process
7 #from complex string which COULD be converted to categories
8 # → to easily processable integers
```

```

9
10 df['date'] = pd.to_datetime(df['date'])
11 df['date'] = df['date'].dt.strftime('%Y.%m.%d')
12 df['year'] = pd.DatetimeIndex(df['date']).year
13 df['month'] = pd.DatetimeIndex(df['date']).month
14 df['day'] = pd.DatetimeIndex(df['date']).day
15
16 print(df.info())

```

```

<class 'pandas.core.frame.DataFrame'>
RangeIndex: 151164 entries, 0 to 151163
Data columns (total 22 columns):
#   Column                                Non-Null Count  Dtype
---  -
0   Unnamed: 0                            151164 non-null  int64
1   iso_code                              151164 non-null  object
2   location                              151164 non-null  object
3   date                                  151164 non-null  object
4   total_cases                           151164 non-null  float64
5   total_deaths                          151164 non-null  int64
6   total_cases_per_million               151164 non-null  float64
7   total_deaths_per_million              151164 non-null  float64
8   total_tests                           151164 non-null  float64
9   total_tests_per_thousand              151164 non-null  float64
10  positive_rate                         151164 non-null  float64
11  people_vaccinated                     151164 non-null  float64
12  people_fully_vaccinated                151164 non-null  float64
13  people_fully_vaccinated_per_hundred    151164 non-null  float64
14  population                             151164 non-null  int64
15  population_density                     151164 non-null  float64
16  median_age                             151164 non-null  float64
17  gdp_per_capita                         151164 non-null  float64
18  human_development_index                151164 non-null  float64
19  year                                   151164 non-null  int64
20  month                                  151164 non-null  int64
21  day                                    151164 non-null  int64
dtypes: float64(13), int64(6), object(3)
memory usage: 25.4+ MB
None

```

```

1 #Defining our X and y to predict country of origin
2 #taking in X
3 #total cases per million
4 #total deaths per million
5 #positive rate
6 #people fully vaccinated per hundred
7 #year
8 #month
9
10 X = df.iloc[:, [6, 7, 10, 13, 19, 20]]
11
12 #taking in y
13 #location to predict
14 y = df.iloc[:, [2]]

```

```

15
16 #train test split , 60:40
17
18 X_train, X_test, y_train, y_test = train_test_split(X, y, test_size=0.4)

1 #Decision Tree Algorithm to predict which country the given data is from
2 #provided an unlabelled dataset.
3
4 country_cdvpyml_clf = DecisionTreeClassifier()
5
6 country_cdvpyml_clf.fit(X_train, y_train)

    DecisionTreeClassifier()

1 y_pred = country_cdvpyml_clf.predict(X_test)
2
3 treeclf_acc = accuracy_score(y_test,y_pred)
4 treeclf_balacc = balanced_accuracy_score(y_test, y_pred)
5 treeclf_prec = precision_score(y_test, y_pred, average='macro')
6 treeclf_jacc = jaccard_score(y_test, y_pred, average='macro')
7
8 print("Decision Tree \naccuracy = ", treeclf_acc, "\nbalanced accuracy = ", treeclf_balac

    Decision Tree
    accuracy =  0.9743326828300202
    balanced accuracy =  0.9744314682819577
    precision =  0.9744823265342877
    jaccard =  0.9508690618074491

1 #Defining our X and y to predict how well vaccinations will go
2 #taking in X
3 #total cases per million
4 #total deaths per million
5 #positive rate
6 #year
7 #month
8
9 X = df.iloc[:, [6, 7, 10, 19, 20]]
10
11 #taking in y
12 #location to predict
13 y = df.iloc[:, [13]]
14
15 #train test split , 60:40
16
17 X_train, X_test, y_train, y_test = train_test_split(X, y, test_size=0.4)

1 scaler = StandardScaler()
2 X_train_scaled = scaler.fit_transform(X_train)

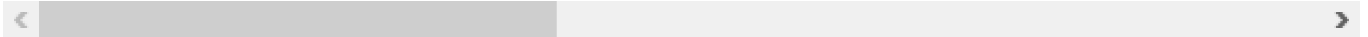
```

```
3 X_test_scaled = scaler.transform(X_test)
```

```
1 #Decision Tree and Random Forest Regression to predict how well vaccinations will go
2 #provided a date and some data at that date
3
4 country_cdpym_dtreg = DecisionTreeRegressor()
5 country_cdpym_rfreg = RandomForestRegressor()
6
7 country_cdpym_dtreg.fit(X_train_scaled, y_train)
8 country_cdpym_rfreg.fit(X_train_scaled, y_train)
```

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

```
RandomForestRegressor()
```



```
1 y_pred = country_cdpym_dtreg.predict(X_test_scaled)
2
3 tree_mse = mean_squared_error(y_test, y_pred)
4 tree_mae = mean_absolute_error(y_test, y_pred)
5
6 y_pred = country_cdpym_rfreg.predict(X_test_scaled)
7
8 rf_mse = mean_squared_error(y_test, y_pred)
9 rf_mae = mean_absolute_error(y_test, y_pred)
10
11 print("Decision Tree \nmse = ",tree_mse,"\nmae = ",tree_mae,"\nrmse = ", sqrt(tree_mse))
12 print("\nRandom Forest \nmse = ",rf_mse," \nmae = ",rf_mae,"\nrmse = ", sqrt(rf_mse))
```

```
Decision Tree
```

```
mse = 6.104453299225523
mae = 0.32586605787972966
rmse = 2.470719186638887
```

```
Random Forest
```

```
mse = 2.9265676986387823
mae = 0.42199349009762854
rmse = 1.7107213971417972
```

```
1 #-----
```

```
1 #Code by Harish
2
3 # creating a copy to perform various opeations
4 df=df_income.copy()
```

```
1 # changing the date format to the int form string to prcoess it easily
2 df['date'] = pd.to_datetime(df['date'])
3 df['date'] = df['date'].dt.strftime('%Y.%m.%d')
```

```

4 df['year'] = pd.DatetimeIndex(df['date']).year
5 df['month'] = pd.DatetimeIndex(df['date']).month
6 df['day'] = pd.DatetimeIndex(df['date']).day
7
8 print(df.info())
<class 'pandas.core.frame.DataFrame'>
RangeIndex: 3864 entries, 0 to 3863
Data columns (total 15 columns):
#   Column                                Non-Null Count  Dtype
---  -
0   Unnamed: 0                            3864 non-null   int64
1   iso_code                              3864 non-null   object
2   location                              3864 non-null   object
3   date                                  3864 non-null   object
4   total_cases                           3864 non-null   int64
5   total_deaths                          3864 non-null   int64
6   total_cases_per_million                3864 non-null   float64
7   total_deaths_per_million               3864 non-null   float64
8   people_vaccinated                      3864 non-null   int64
9   people_fully_vaccinated                3864 non-null   int64
10  people_fully_vaccinated_per_hundred    3864 non-null   float64
11  population                              3864 non-null   int64
12  year                                    3864 non-null   int64
13  month                                    3864 non-null   int64
14  day                                     3864 non-null   int64
dtypes: float64(3), int64(9), object(3)
memory usage: 452.9+ KB
None

```

```

1 # X (labels) -->
2 # x1 = total cases per million
3 # x2 = total deaths per million
4 # x3 = people fully vaccinated per hundred
5 # x4 = year
6 # x5 = month
7
8 X = df.iloc[:, [6, 7, 10, 12, 13]]
9
10 # predicting y -> regions affected
11 y = df.iloc[:, [2]]
12
13 # splitting the train(60) test(40) data
14 X_train, X_test, y_train, y_test = train_test_split(X, y, test_size=0.4)

```

```

1 # Decision Tree Algorithm to predict which country the given data is from
2 country_cdvpyml_clf = DecisionTreeClassifier()
3 country_cdvpyml_clf.fit(X_train, y_train)

DecisionTreeClassifier()

1 y_pred = country_cdvpyml_clf.predict(X_test)

```

```

2 print(y_pred)
3
4 acc = accuracy_score(y_test,y_pred)
5 bal_acc = balanced_accuracy_score(y_test, y_pred)
6 pre = precision_score(y_test, y_pred, average='macro')
7 jacc = jaccard_score(y_test, y_pred, average='macro')
8
9 print("Decision Tree \naccuracy = ", acc, "\nbalanced accuracy = ", bal_acc, "\nprecisior
    ['High income' 'High income' 'High income' ... 'High income'
     'Lower middle income' 'High income']
Decision Tree
accuracy = 0.9844760672703752
balanced accuracy = 0.9847029487974239
precision = 0.9846308971009321
jaccard = 0.9698604184896511

```

```

1 # X (labels) -->
2 # x1 = total cases per million
3 # x2 = total deaths per million
4 # x4 = year
5 # x5 = month
6
7 X = df.iloc[:, [6, 7, 12, 13]]
8
9 # y = people fully vaccinated per hundred (predcited)
10 y = df.iloc[:, [10]]
11
12 # splitting the train(60) test(40) data
13 X_train, X_test, y_train, y_test = train_test_split(X, y, test_size=0.4)

```

```

1 scaler = StandardScaler()
2 X_train_scaled = scaler.fit_transform(X_train)
3 X_test_scaled = scaler.transform(X_test)

```

```

1 # Desicion tree algorithm to predict the y (no. of people gets vaccinated)
2 region_dt = DecisionTreeRegressor()
3 region_dt.fit(X_train_scaled, y_train)
4
5 y_pred = region_dt.predict(X_test_scaled)
6 print(y_pred)
7
8 dt_mse = mean_squared_error(y_test, y_pred)
9 dt_mae = mean_absolute_error(y_test, y_pred)

```

```
[7.656e+01 7.839e+01 2.946e+01 ... 1.820e+00 2.000e-02 5.930e+00]
```

```

1 # Random Forest algorithm to predict the y (no. of people gets vaccinated)
2 region_rf = RandomForestRegressor()
3 region_rf.fit(X_train_scaled, y_train)

```

```

4
5 y_pred = reigion_rf.predict(X_test_scaled)
6 print(y_pred)
7
8 rf_mse = mean_squared_error(y_test, y_pred)
9 rf_mae = mean_absolute_error(y_test, y_pred)

/usr/local/lib/python3.7/dist-packages/ipykernel_launcher.py:3: DataConversionWarning: /
  This is separate from the ipykernel package so we can avoid doing imports until
[7.65607e+01 7.84550e+01 2.98320e+01 ... 1.80680e+00 1.75000e-02
 6.14680e+00]

```

```

1 # comparing the errors in algorithm
2
3 print("Decision Tree \nmse = ",dt_mse,"\nmae = ",dt_mae,"\nrmse = ", sqrt(dt_mse))
4 print("\nRandom Forest \nmse = ",rf_mse," \nmae = ",rf_mae,"\nrmse = ", sqrt(rf_mse))

```

```

Decision Tree
mse = 0.0303601552393273
mae = 0.05817593790426913
rmse = 0.1742416575888995

```

```

Random Forest
mse = 0.06462492855109962
mae = 0.05394320827943078
rmse = 0.2542143358489045

```

```
1 #-----
```

```

1 # Code by Siddh
2
3 # creating a copy to perform various opeations
4 df=df_region.copy()

```

```

1 # changing the date format to the int form string to prcoess it easily
2 df['date'] = pd.to_datetime(df['date'])
3 df['date'] = df['date'].dt.strftime('%Y.%m.%d')
4 df['year'] = pd.DatetimeIndex(df['date']).year
5 df['month'] = pd.DatetimeIndex(df['date']).month
6 df['day'] = pd.DatetimeIndex(df['date']).day
7
8 print(df.info())

```

```

<class 'pandas.core.frame.DataFrame'>
RangeIndex: 7734 entries, 0 to 7733
Data columns (total 15 columns):
#   Column                                Non-Null Count  Dtype
---  -
0   Unnamed: 0                            7734 non-null   int64
1   iso_code                              7734 non-null   object

```

```

2 location 7734 non-null object
3 date 7734 non-null object
4 total_cases 7734 non-null int64
5 total_deaths 7734 non-null int64
6 total_cases_per_million 7734 non-null float64
7 total_deaths_per_million 7734 non-null float64
8 people_vaccinated 7734 non-null int64
9 people_fully_vaccinated 7734 non-null int64
10 people_fully_vaccinated_per_hundred 7734 non-null float64
11 population 7734 non-null int64
12 year 7734 non-null int64
13 month 7734 non-null int64
14 day 7734 non-null int64
dtypes: float64(3), int64(9), object(3)
memory usage: 906.5+ KB
None

```

```

1 # X (labels) -->
2 # x1 = total cases per million
3 # x2 = total deaths per million
4 # x3 = people fully vaccinated per hundred
5 # x4 = year
6 # x5 = month
7
8 X = df.iloc[:, [6, 7, 10, 12, 13]]
9
10 # predicting y -> regions affected
11 y = df.iloc[:, [2]]
12
13 # splitting the train(60) test(40) data
14 X_train, X_test, y_train, y_test = train_test_split(X, y, test_size=0.4)

1 # Decision Tree Algorithm to predict which country the given data is from
2 country_cdvpyml_clf = DecisionTreeClassifier()
3 country_cdvpyml_clf.fit(X_train, y_train)

DecisionTreeClassifier()

1 y_pred = country_cdvpyml_clf.predict(X_test)
2 print(y_pred)
3
4 acc = accuracy_score(y_test, y_pred)
5 bal_acc = balanced_accuracy_score(y_test, y_pred)
6 pre = precision_score(y_test, y_pred, average='macro')
7 jacc = jaccard_score(y_test, y_pred, average='macro')
8
9 print("Decision Tree \naccuracy = ", acc, "\nbalanced accuracy = ", bal_acc, "\nprecisior

['Europe' 'Asia' 'Africa' ... 'South America' 'World' 'Asia']
Decision Tree
accuracy = 0.9091790562378798

```



```

balanced accuracy = 0.9096047809120531
precision = 0.9098267920680132
jaccard = 0.8379172876900025

```

```

1 # X (labels) -->
2 # x1 = total cases per million
3 # x2 = total deaths per million
4 # x4 = year
5 # x5 = month
6
7 X = df.iloc[:, [6, 7, 12, 13]]
8
9 # y = people fully vaccinated per hundred (predcited)
10 y = df.iloc[:, [10]]
11
12 # splitting the train(60) test(40) data
13 X_train, X_test, y_train, y_test = train_test_split(X, y, test_size=0.4)

```

```

1 scaler = StandardScaler()
2 X_train_scaled = scaler.fit_transform(X_train)
3 X_test_scaled = scaler.transform(X_test)

```

```

1 # Desicion tree algorithm to predict the y (no. of people gets vaccinated)
2 reigion_dt = DecisionTreeRegressor()
3 reigion_dt.fit(X_train_scaled, y_train)
4
5 y_pred = reigion_dt.predict(X_test_scaled)
6 print(y_pred)
7
8 dt_mse = mean_squared_error(y_test, y_pred)
9 dt_mae = mean_absolute_error(y_test, y_pred)

```

```
[15.73  0.06 44.22 ... 41.32 13.15 48.05]
```

```

1 # Random Forest algorithm to predict the y (no. of people gets vaccinated)
2 reigion_rf = RandomForestRegressor()
3 reigion_rf.fit(X_train_scaled, y_train)
4
5 y_pred = reigion_rf.predict(X_test_scaled)
6 print(y_pred)
7
8 rf_mse = mean_squared_error(y_test, y_pred)
9 rf_mae = mean_absolute_error(y_test, y_pred)

```

```

/usr/local/lib/python3.7/dist-packages/ipykernel_launcher.py:3: DataConversionWarning: /
This is separate from the ipykernel package so we can avoid doing imports until
[15.7029  0.0565 44.5292 ... 42.7917 13.0369 51.1516]

```



```
1 # comparing the errors in algorithm
2
3 print("Decision Tree \nmse = ",dt_mse,"\nmae = ",dt_mae,"\nrmse = ", sqrt(dt_mse))
4 print("\nRandom Forest \nmse = ",rf_mse," \nmae = ",rf_mae,"\nrmse = ", sqrt(rf_mse))
```

Decision Tree

mse = 1.8603172268907564

mae = 0.2360730446024564

rmse = 1.3639344657610044

Random Forest

mse = 1.1992298690723981

mae = 0.2672010019392378

rmse = 1.0950935435260305

[Colab paid products](#) - [Cancel contracts here](#)

✓ 0s completed at 11:46 PM

