## AIM: Implementation of Bagging Algorithm: Decision Tree, Random Forest

**THEORY:**

### 1) Decision Tree:

A decision tree is a flowchart-like structure in which each internal node represents a test on a feature (e.g. whether a coin flip comes up heads or tails) , each leaf node represents a class label (decision taken after computing all features) and branches represent conjunctions of features that lead to those class labels. The paths from root to leaf represent classification rules. Below diagram illustrate the basic flow of decision tree for decision making with labels (Rain(Yes), No Rain(No)).

Decision tree is one of the predictive modelling approaches used in statistics, data mining and machine learning.

Decision trees are constructed via an algorithmic approach that identifies ways to split a data set based on different conditions. It is one of the most widely used and practical methods for supervised learning. Decision Trees are a nonparametric supervised learning method used for both classification and regression tasks.

Tree models where the target variable can take a discrete set of values are called classification trees. Decision trees where the target variable can take continuous values (typically real numbers) are called regression trees.

### 2) Random Forest:

Random forest is an ensemble machine learning algorithm.

It is perhaps the most popular and widely used machine learning algorithm given its good or excellent performance across a wide range of classification and regression predictive modeling problems.

It is also easy to use given that it has few key hyperparameters and sensible heuristics for configuring these hyperparameters.

It is an extension of bootstrap aggregation (bagging) of decision trees and can be used for classification and regression problems.

In bagging, a number of decision trees are created where each tree is created from a different bootstrap sample of the training dataset. A bootstrap sample is a sample of the training dataset where a sample may appear more than once in the sample, referred to as sampling with replacement.

Bagging is an effective ensemble algorithm as each decision tree is fit on a slightly different training dataset, and in turn, has a slightly different performance. Unlike normal decision tree models, such as classification and regression trees (CART), trees used in the ensemble are unpruned, making them slightly overfit to the training dataset. This is desirable as it helps to make each tree more different and have less correlated predictions or prediction errors.

Predictions from the trees are averaged across all decision trees resulting in better performance than any single tree in the model.

### 1) <u>IMPORTING LIBRARIES:</u>

```python
import numpy as np # linear algebra
import pandas as pd # data processing, CSV file I/O (e.g. pd.read_csv)
import seaborn as sns
import matplotlib.pyplot as plt
```

### 2) <u>READING DATASET [TRAINING & TESTING]:</u>

```python
train_df = pd.read_csv("Disease Training.csv")
test_df = pd.read_csv("Disease Testing.csv")
train_df.head()
```

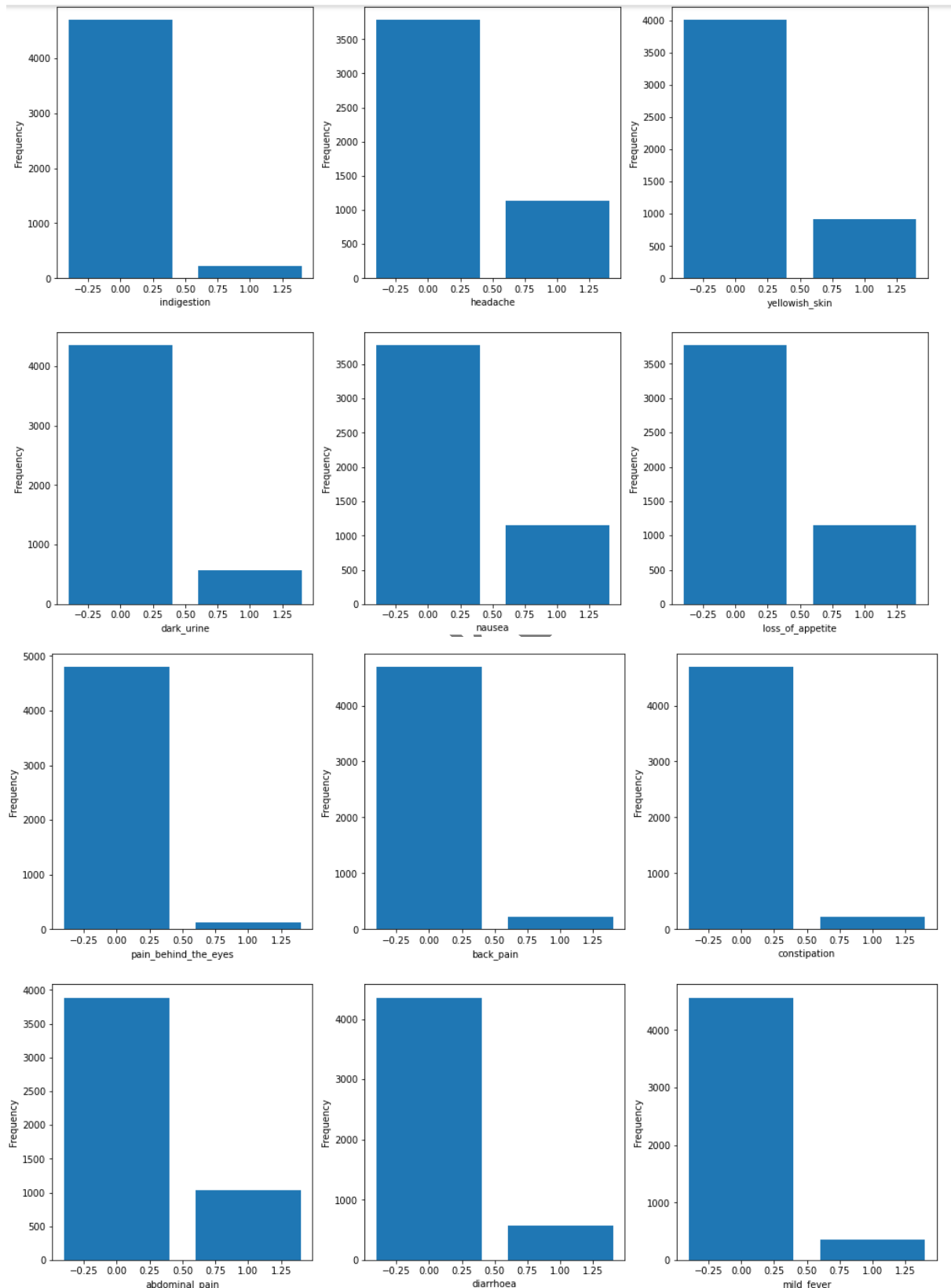| | itching | skin_rash | nodal_skin_eruptions | continuous_sneezing | shivering | chills | joint_pain | stomach_pain | acidity | ulcers_on_tongue | ... | scurring | skin_peeling | silver_like_dusting | small_dents_in_nails |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| 0 | 1 | 1 | 1 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | ... | 0 | 0 | 0 | 0 |
| 1 | 0 | 1 | 1 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | ... | 0 | 0 | 0 | 0 |
| 2 | 1 | 0 | 1 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | ... | 0 | 0 | 0 | 0 |
| 3 | 1 | 1 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | ... | 0 | 0 | 0 | 0 |
| 4 | 1 | 1 | 1 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | ... | 0 | 0 | 0 | 0 |

5 rows × 134 columns

### 3) <u>DATA CLEANING:</u>

```python
# drop unnamed feature from train data
train_df.drop("Unnamed: 133", axis = 1, inplace = True)
# train_df["Unnamed: 133"]  # it's not here anymore
```
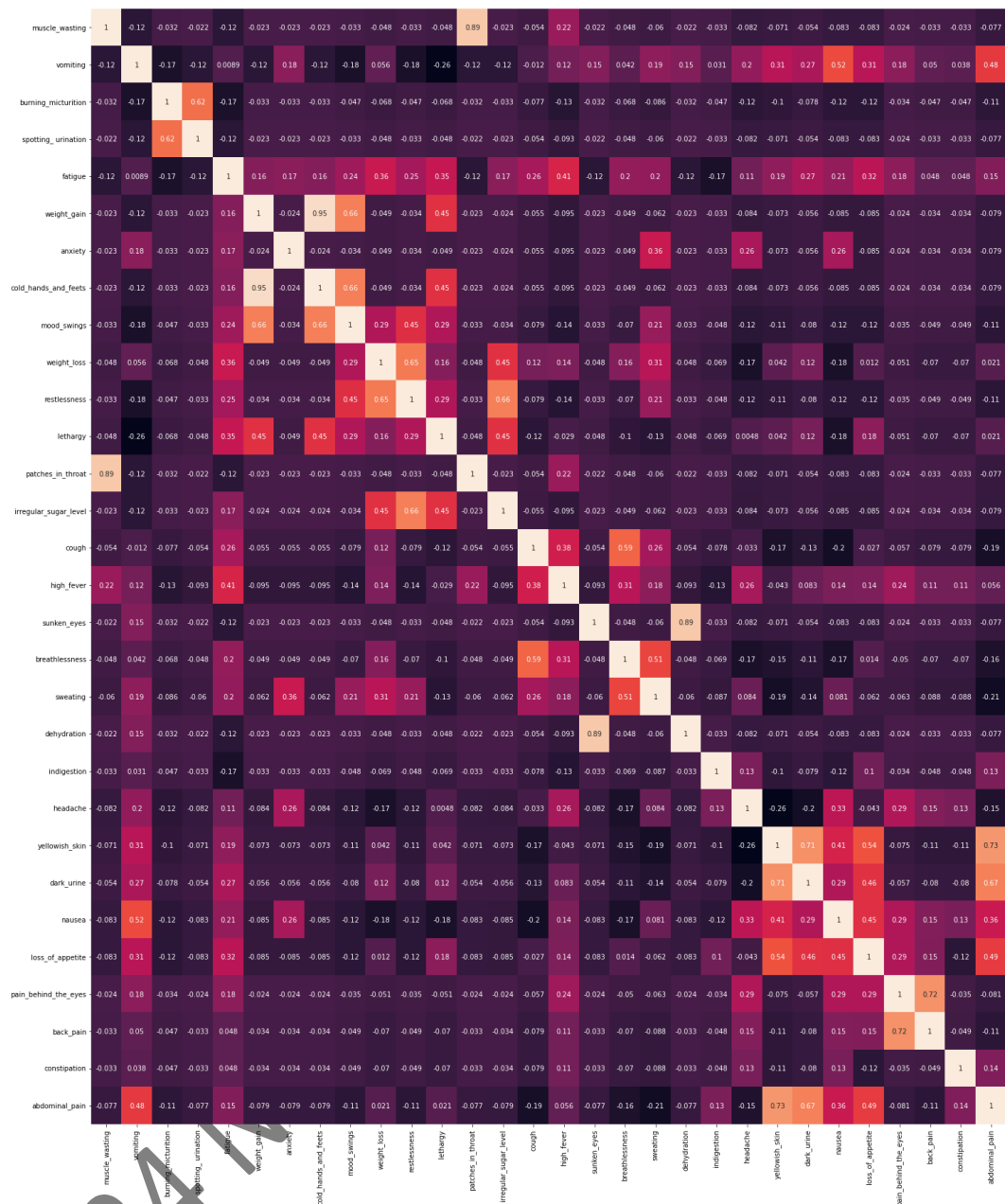
### 4) <u>FEATURES VISULALIZATION:</u>

```python
# lets visualize some of features
features = ['indigestion', 'headache', 'yellowish_skin', 'dark_urine', 'nausea',
        'loss_of_appetite', 'pain_behind_the_eyes', 'back_pain', 'constipation',
        'abdominal_pain', 'diarrhoea', 'mild_fever']

plt.figure(figsize = (17,25))
for i, feature in enumerate(features):
    plt.subplot(4,3,i+1)
    plt.bar(train_df[feature].value_counts().index.to_numpy(), train_df[feature].value_counts().values)
    plt.xlabel(feature)
    plt.ylabel("Frequency")
plt.show()
```

**FYMCA-B**
**AL/ML**

**SEM-II**
**PRACTICAL NO: 09**

**DATE: 27/06/2022**
**ROLL NO: 24**

```python
# linear relationships between some of features using correlation heatmap: for example which symptoms occur together?
df_corr = train_df.iloc[:, 10:40]
plt.figure(figsize = (30, 30))
sns.heatmap(df_corr.corr(), annot = True)
plt.show()
```

**FYMCA-B**
**AL/ML**

**SEM-II**
**PRACTICAL NO: 09**

**DATE: 27/06/2022**
**ROLL NO: 24**

### 5) BUILDING RANDOM FOREST MODEL:

```python
from sklearn.ensemble import RandomForestClassifier
x_train, y_train = train_df.loc[:,train_df.columns != "prognosis"], train_df.loc[:,"prognosis"]
x_test, y_test = test_df.loc[:,train_df.columns != "prognosis"], test_df.loc[:,"prognosis"]
rfc = RandomForestClassifier(random_state = 42, n_estimators = 100)
rfc.fit(x_train, y_train)
rfc.predict(x_test)
```

```
array(['Fungal infection', 'Allergy', 'GERD', 'Chronic cholestasis',
       'Drug Reaction', 'Peptic ulcer diseae', 'AIDS', 'Diabetes ',
       'Gastroenteritis', 'Bronchial Asthma', 'Hypertension ', 'Migraine',
       'Cervical spondylosis', 'Paralysis (brain hemorrhage)', 'Jaundice',
       'Malaria', 'Chicken pox', 'Dengue', 'Typhoid', 'hepatitis A',
       'Hepatitis B', 'Hepatitis C', 'Hepatitis D', 'Hepatitis E',
       'Alcoholic hepatitis', 'Tuberculosis', 'Common Cold', 'Pneumonia',
       'Dimorphic hemmorhoids(piles)', 'Heart attack', 'Varicose veins',
       'Hypothyroidism', 'Hyperthyroidism', 'Hypoglycemia',
       'Osteoarthristis', 'Arthritis',
       '(vertigo) Paroymsal  Positional Vertigo', 'Acne',
       'Urinary tract infection', 'Psoriasis', 'Impetigo', 'Impetigo'],
      dtype=object)
```

### 6) CHECKING RANDOM FOREST SCORE:

```python
rfc.score(x_test, y_test)
```

```
0.9761904761904762
```

### 7) BUILDING BAGGING CLASSIFIER:

```python
from sklearn.metrics import confusion_matrix, accuracy_score, classification_report
from sklearn.ensemble import BaggingClassifier
from sklearn.tree import DecisionTreeClassifier
bagging_clf = BaggingClassifier(base_estimator=tree, n_estimators=1500,
random_state=42)
bagging_clf.fit(x_train, y_train)
y_test_pred = bagging_clf.predict(x_test)
y_train_pred = bagging_clf.predict(x_train)
print("TRAINIG RESULTS: \n===============================")
clf_report = pd.DataFrame(classification_report(y_train, y_train_pred,
output_dict=True))
print(f"CONFUSION MATRIX:\n{confusion_matrix(y_train, y_train_pred)}")
print(f"ACCURACY SCORE:\n{accuracy_score(y_train, y_train_pred):.4f}")
print(f"CLASSIFICATION REPORT:\n{clf_report}")
print("TESTING RESULTS: \n===============================")
clf_report = pd.DataFrame(classification_report(y_test, y_test_pred,
output_dict=True))
print(f"CONFUSION MATRIX:\n{confusion_matrix(y_test, y_test_pred)}")
print(f"ACCURACY SCORE:\n{accuracy_score(y_test, y_test_pred):.4f}")
print(f"CLASSIFICATION REPORT:\n{clf_report}")
```

```
TRAINIG RESULTS:
================================
CONFUSION MATRIX:
[[120   0   0 ...   0   0   0]
 [  0 120   0 ...   0   0   0]
 [  0   0 120 ...   0   0   0]
 ...
 [  0   0   0 ... 120   0   0]
 [  0   0   0 ...   0 120   0]
 [  0   0   0 ...   0   0 120]]
ACCURACY SCORE:
1.0000
CLASSIFICATION REPORT:
           (vertigo) Paroymsal Positional Vertigo   AIDS   Acne  \
precision                                      1.0    1.0    1.0
recall                                         1.0    1.0    1.0
f1-score                                       1.0    1.0    1.0
support                                      120.0  120.0  120.0

           Alcoholic hepatitis  Allergy  Arthritis  Bronchial Asthma  \
precision                  1.0      1.0        1.0               1.0
recall                     1.0      1.0        1.0               1.0
f1-score                   1.0      1.0        1.0               1.0
support                  120.0    120.0      120.0             120.0

           Cervical spondylosis  Chicken pox  Chronic cholestasis  ...  \
precision                   1.0          1.0                  1.0  ...
recall                      1.0          1.0                  1.0  ...
f1-score                    1.0          1.0                  1.0  ...
support                   120.0        120.0                120.0  ...

           Pneumonia  Psoriasis  Tuberculosis  Typhoid  \
precision        1.0        1.0           1.0      1.0
recall           1.0        1.0           1.0      1.0
f1-score         1.0        1.0           1.0      1.0
support        120.0      120.0         120.0    120.0

           Urinary tract infection  Varicose veins  hepatitis A  accuracy  \
precision                      1.0             1.0          1.0       1.0
recall                         1.0             1.0          1.0       1.0
f1-score                       1.0             1.0          1.0       1.0
support                      120.0           120.0        120.0       1.0

           macro avg  weighted avg
precision        1.0           1.0
recall           1.0           1.0
f1-score         1.0           1.0
support       4920.0        4920.0

[4 rows x 44 columns]
```

```
TESTING RESULTS:
================================
CONFUSION MATRIX:
[[1 0 0 ... 0 0 0]
 [0 1 0 ... 0 0 0]
 [0 0 1 ... 0 0 0]
 ...
 [0 0 0 ... 1 0 0]
 [0 0 0 ... 0 1 0]
 [0 0 0 ... 0 0 1]]
ACCURACY SCORE:
0.9762
CLASSIFICATION REPORT:
           (vertigo) Paroymsal Positional Vertigo  AIDS  Acne  \
precision                                     1.0   1.0   1.0
recall                                        1.0   1.0   1.0
f1-score                                      1.0   1.0   1.0
support                                       1.0   1.0   1.0

           Alcoholic hepatitis  Allergy  Arthritis  Bronchial Asthma  \
precision                  1.0      1.0        1.0               1.0
recall                     1.0      1.0        1.0               1.0
f1-score                   1.0      1.0        1.0               1.0
support                    1.0      1.0        1.0               1.0

           Cervical spondylosis  Chicken pox  Chronic cholestasis  ...  \
precision                   1.0     0.500000                  1.0  ...
recall                      1.0     1.000000                  1.0  ...
f1-score                    1.0     0.666667                  1.0  ...
support                     1.0     1.000000                  1.0  ...

           Pneumonia  Psoriasis  Tuberculosis  Typhoid  \
precision        1.0        1.0           1.0      1.0
recall           1.0        1.0           1.0      1.0
f1-score         1.0        1.0           1.0      1.0
support          1.0        1.0           1.0      1.0

           Urinary tract infection  Varicose veins  hepatitis A  accuracy  \
precision                      1.0             1.0          1.0   0.97619
recall                         1.0             1.0          1.0   0.97619
f1-score                       1.0             1.0          1.0   0.97619
support                        1.0             1.0          1.0   0.97619

           macro avg  weighted avg
precision   0.987805      0.988095
recall      0.987805      0.976190
f1-score    0.983740      0.976190
support    42.000000     42.000000

[4 rows x 44 columns]
```
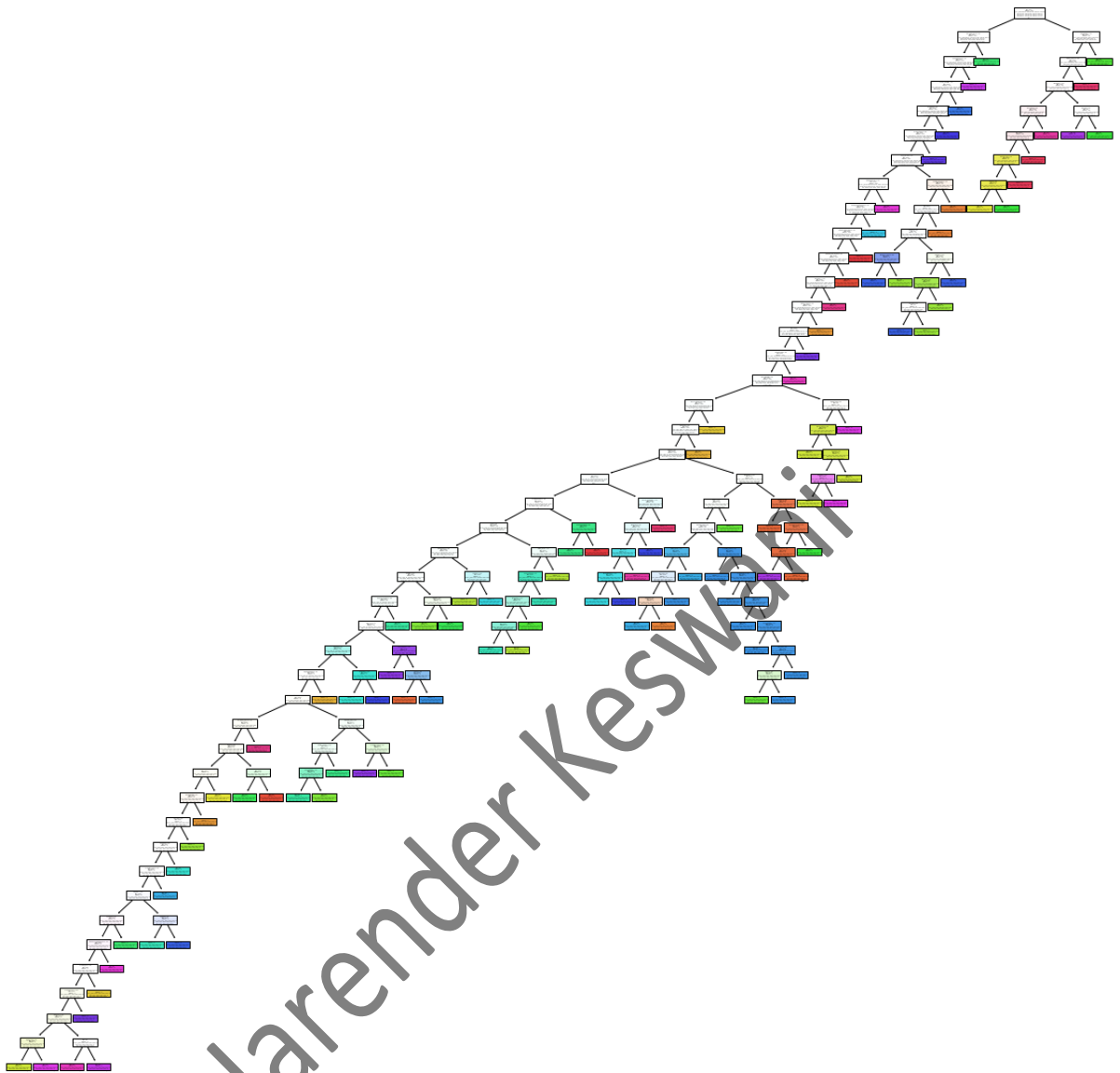
## 8) PLOTTING TREE:

```
!pip install dtreeviz
from dtreeviz.trees import dtreeviz # will be used for tree visualization
from matplotlib import pyplot as plt
from sklearn import tree
plt.figure(figsize=(20,20))
X = pd.DataFrame(train_df, columns=train_df.columns)
_ = tree.plot_tree(rfc.estimators_[0], feature_names=X.columns, filled=True)
```

**CONCLUSION:**

From this practical, I have learned and implemented the random forest algorithm in python.