| | |
|---|---|
| Course Code: CT4101 (Machine Learning) | Assignment: 2 |
| Name: *Seshadri Sundarrajan*    Course Code: 1CSD1 | Student ID: 19230107 |

## RANDOM FOREST FROM SCRATCH - REPORT                        (DATASET:HAZZLENUTS)

- Implemented Random Forest algorithm encompassing Bootstrap sampling and random feature subspace. Each classifier(tree) of the algorithm was built on C4.5 algorithm.
- Random forest is an ensemble technique and known for better accuracy as it combines output from several weak learners to form base classifier. Moreover, once Decision tree classifier is implemented it requires minimal effort to build bootstrap aggregation methodology on top of it to get bagged model.

## PROGRAMMING LANGUAGE AND LIBRARIES USED

- **Python**
    - Primarily used *Pandas* Data frame to store tree information (Meta data).
    - For computations made use of Numpy Arrays.
    - Matplotlib is used for visualization ROC curves.
    - Networkx is used for decision tree visualization.

## DATA PRE-PROCESSING (DONE IN EXCEL)

- The provided dataset had columns representing instances. To feed into Pandas data frame and for ease of interpreting *transposed* the data to get instances along the row.
- 'Sample id' column has no impact on our classifier as it's just a key, so dropped it.

## DESIGN LOGIC

### RANDOM FOREST

Algorithm is built upon 2 random processes

1. *Bootstrap sampling*
    - At each tree, input data is split to n-samples with repetition.
2. *Feature Subspace*
    - At each tree, features are subset randomly

**Parameters Considered for design**: These are part of constructor (should be passed while creating object)

| Parameters | Description |
|---|---|
| Total_Trees | Number of decision trees to be constructed |
| Tree Depth | Maximum depth of a decision tree |
| Bootstrap size | Bootstrap sample size per tree |
| Features_per_tree | Random number of features on which a tree is built |
| Random_state | Random seed |

**Training.**

Each individual tree is constructed using C4.5 algorithm implemented separately. As part of C4.5 design we get meta information on each tree and we store it.

**Testing**

For each test data, we get the output from each tree. Finally, output of a single test data will be aggregation of all the collected tree outputs. (Mode).

**Convergence**: All trees are trained and individual tree convergence depends on Decision Tree convergence.

**Methods Available:**

| Name | Parameters | Return | Description |
|---|---|---|---|
| Fit() | - | - | Starts building model. |
| Predict() | Array (x_test) | Predictions(array) | Predicts output class for each element in the input array from the built model |
| Predict_prob() | Array (x_test) | Output class probabilities (array) | Returns class probabilities (Mean probability across all tress for each class) - To get ROC Curves |
| Write_meta_data() | File_name | Excel workbook | Writes each tree meta info to a sheet. |

**NOTE**: Private method named '**Predict_tree_op**' is called inside 'Predict' to get output class from an individual tree for given test data. So, can we can perform aggregation on top of all tree outputs.

## C4.5

Algorithm is built entirely to support *continuous* features and uses *entropy* as impurity measure.

**Binning of continuous features**

- For given feature, set of thresholds are chosen for evaluation by taking mean between successive values arranged in ascending order. (For n values we get n-1 threshold).
- For each threshold we calculate its entropy and overall entropy based on class proportions.
- Finally, best threshold that has minimum overall entropy is chosen to get homogeneous splits with high information gain. We split data into 2 partitions namely (<=Threshold & >Threshold).

**Meta data driven Approach**

1. We select the best feature that gives us homogenous split and start navigating down.
2. After the first split we store the corresponding information into our meta data frame.
3. Then next corresponding split takes place by accessing the subset of data that branched out and selecting the necessary feature and its threshold from the meta info.
4. Thus, each split (except first split) is dependent on the meta data of previous split.
5. This process keeps reoccurring until convergence.

| | | | | META DATA FRAME STRUCTURE | | | | | | | |
|---|---|---|---|---|---|---|---|---|---|---|---|

| node_number | parent_node | split_feature | split_category | split_condition | overall_entropy | leaf_f | next_split_features | c_cornuta | _american | c_avellana |
|---|---|---|---|---|---|---|---|---|---|---|
| 1 | | shell_radius | | [5.0565] | 0.824687159 | N | ['hardness', 'carbohydrate', 'length', 'compactness', 'mass'] | 0.34 | 0.33 | 0.33 |
| 1_1 | 1 | shell_radius | <= 5.0565 | ['shell_radius <= 5.0565'] | 0.383942147 | N | ['hardness', 'carbohydrate', 'length', 'compactness', 'mass'] | 0.015385 | 0.476923 | 0.507692 |
| 1_2 | 1 | shell_radius | > 5.0565 | ['shell_radius > 5.0565'] | 0.114285714 | N | ['hardness', 'carbohydrate', 'length', 'compactness', 'mass'] | 0.942857 | 0.057143 | 0 |
| 1_1_1 | 1_1 | mass | <= 1467.9 | ['shell_radius <= 5.0565', 'mass <= 1467.9'] | 0.481178121 | N | ['hardness', 'carbohydrate', 'length', 'compactness'] | 0.026316 | 0.105263 | 0.868421 |
| 1_1_2 | 1_1 | mass | > 1467.9 | ['shell_radius <= 5.0565', 'mass > 1467.9'] | 0 | Y | ['hardness', 'carbohydrate', 'length', 'compactness'] | 0 | 1 | 0 |

*'Split condition'* is list of strings (stored cummulatively at each node depending on its branch hierarchy). Ex: - ['shell radius <= 5.012','width<=14.148'], this will drive us to get the subset needed for next split.

*'Next split feature'* is a list with feature names upon which successive splits can be done. They are decided in such a way that they do not repeat the same feature that branched out to the current node.

### Recursive Logic

1. At each non-leaf node, we get the node's split condition and filter the data subset for next split.
2. Then from *'next_split_features'* find best feature and its threshold along with overall entropy.
3. Then we populate child nodes data along with output class probabilities and leaf flag(**'leaf_f'**).
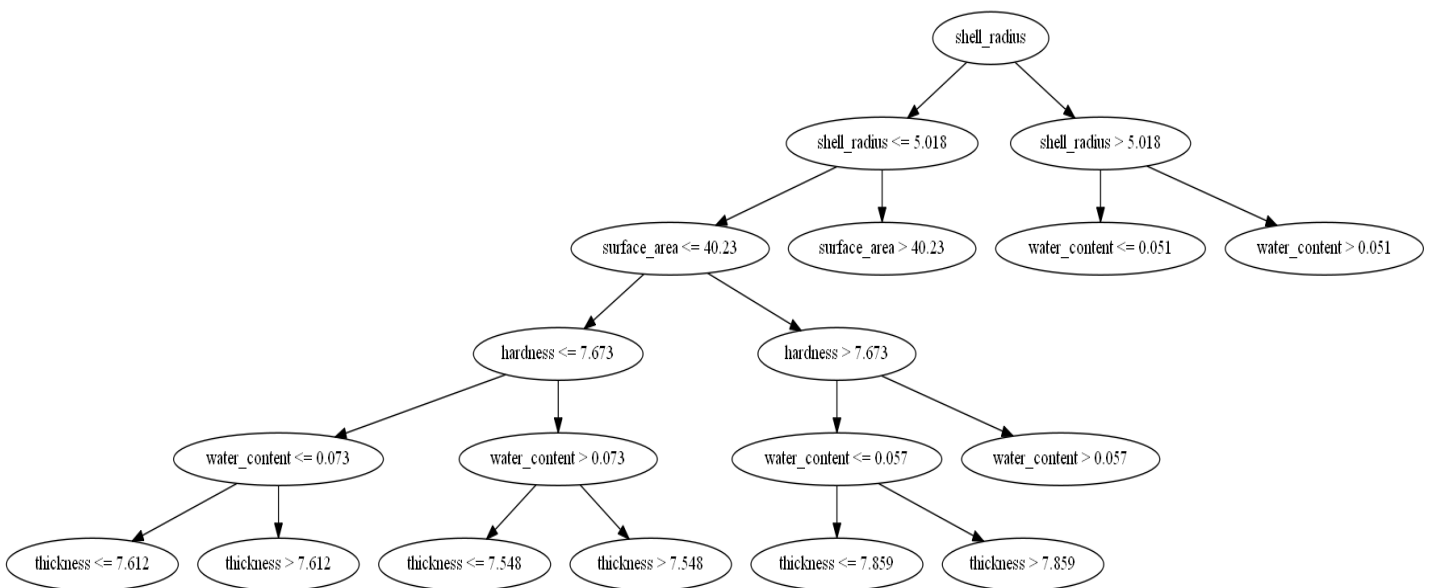
### Convergence

- If a node is pure then we label them as 'leaf' in our meta data, so that no further split occurs.
- On contrary, if our tree depth is reached and yet node is not pure, again we label them as 'leaf'.

**Parameters Considered for design:** Tree depth

**Methods Available:** Fit() [ Starts building the tree and returns the tree meta information as data frame].

**NOTE**: Private methods named *'tree_initial_split'*, *'build_tree'* and *'recursive_tree'* are used to build the tree and return tree information.

## DECISION TREE VISUALIZED FROM META DATA (PACKAGE USED: NETWORKX)



## RANDOM FOREST (SCI-KIT LEARN)

Created model from ML library to replicate the same behaviour as that of implementation done, where test and train data are split in ratio 1:2.

| Parameters | Scratch Implementation | Sklearn library |
|---|---|---|
| Total_Trees | 11 | 11 |
| Tree Depth | 5 | 5 |
| Bootstrap size | 100 | Not provided |
| Features_per_tree | 6 | Not provided |
| Random_state | 10 | 10 |

- Sklearn Model execution **time** was much faster than the implementation. Potential reasons might be caching and parallel processing.
- Initial comparison was done at random state =10 and **confusion matrices** obtained are displayed.

```
Confusion Matrix Random Forest (Implementation from scratch)

             c_cornuta  c_avellana  c_americana
c_cornuta          14           5            0
c_avellana          1          24            0
c_americana         2           0           21

Accuracy = 88.06%
```
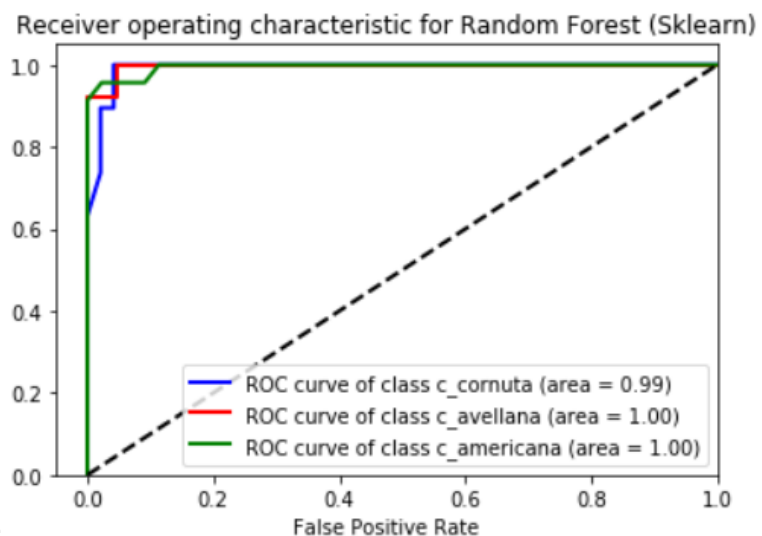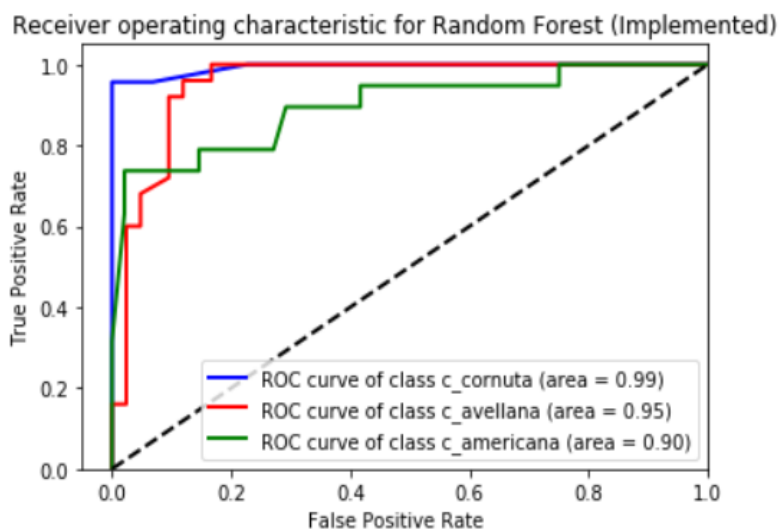
```
Confusion Matrix Random Forest (Implementation from Sklearn)

             c_cornuta  c_avellana  c_americana
c_cornuta          19           0            0
c_avellana          1          23            1
c_americana         1           0           22

Accuracy = 95.52%
```

- Mean accuracy after 10 random splits is *90% for sklearn* random forest and is **88% for implementation** done.
- ROC cures for both the models are generated at random state = 10 and we observe that area under the curve for sklearn random forest is almost 100% for all the 3 classes but in implementation done we have 99% for 'c_cornuta' and 90% for 'c_americana' and 95% for 'c_avellana'. Thus, model implemented tend to misclassify last 2 classes.



- Another important difference between 2 implementations is that there is no interface in sklearn random forest to specify bootstrap sample and feature size per tree, as it is taken care dynamically.
- On comparison to Assignment 1, here haven't removed any correlated features because while generating random feature subset we get a chance of picking correlated feature while its dependent is left out, so our predictions might improve.
- To sum up, scratch implementation could meet ML library performance if it is tuned better with extra parameters and if we calculate bootstrap sample and feature size per tree dynamically from data based on proven methods and taking in to account dynamic programming to compensate time delay. Overall, **Sklearn Random Forrest has better performance**.

## REFERENCES

1. https://www.sebastian-mantey.com/posts/random-forest-algorithm
2. https://books.google.ie/books?id=HExncpjbYroC&printsec=frontcover&source=gbs_ge_summary_r&redir_esc=y#v=onepage&q&f=false
3. https://www.sebastian-mantey.com/posts/decision-tree-algorithm-part-2-entropy
4. https://scikit-learn.org/stable/modules/generated/sklearn.ensemble.RandomForestClassifier.html
5. https://scikit-learn.org/stable/auto_examples/model_selection/plot_roc.html

## APPENDIX – SOURCE CODE

**NOTE**: Runnable file is '**__main__.py** '(the main method file) and it makes the call to algorithm.

## PACKAGE STRUCTURE



## RANDOM_FOREST.PY

```python
1.  import random
2.  import pandas as pd
3.  import numpy as np
4.  import pprint
5.  from operator import itemgetter
6.
7.  # Files for this below 2 imports are attached
8.
9.  from Algorithm.Decision_Tree.C4_5 import * # CART implementation from scratch
10. from Algorithm.Utilities.Calculations import *   # All helper functions involving calculations
11.
12.
13.
14.
15. class random_forest():
16.
```

```python
17.      # At each tree we store features sent and their indexes which will be used for testing
18.      features_idx_subseted = []
19.      feature_subseted = []
20.
21.      # We store Meta info in a dictionary with assumed structure { 'Tree_1' : 'Meta_Data_frame_of_t
   ree1',.....}
22.      meta_dfs = {}
23.      def __init__ (self,data,total_trees,bootstrap_size,feature_per_tree,tree_depth = 5,random_seed
   = 10):
24.
25.          self.feature_per_tree = feature_per_tree
26.          self.bootstrap_size = bootstrap_size
27.          self.total_trees = total_trees
28.          self.data = data
29.          self.features =data.columns.tolist()[0:-1]
30.          self.lable_name = data.columns.tolist()[-1]
31.          self.tree_depth = tree_depth
32.          self.random_seed = random_seed
33.
34.
35.      def fit(self):
36.
37.          '''''Build each tree using Bootstrap sample Random Feature subspace.'''
38.
39.          random.seed(self.random_seed)
40.
41.          for tree in range(self.total_trees):
42.              sample = get_bootstrap(self.data,self.bootstrap_size,self.random_seed)
43.
44.              sub_feature = random.sample(range(len(self.features)), k=self.feature_per_tree)
45.              random_forest.features_idx_subseted.append(sub_feature)
46.
47.              feature_subset = list(itemgetter(*sub_feature)(self.features))
48.              random_forest.feature_subseted.append(feature_subset)
49.
50.              # Call to decision tree
51.
52.              c = DT_classifier(self.tree_depth,sample,feature_subset,self.lable_name,sample[self.la
   ble_name].unique().tolist())
53.
54.              c.fit()
55.
56.
57.              # Appending decision tree meta info of this tree to our random forest meta dictionary
58.              random_forest.meta_dfs[tree] = c.mdf
59.
60.
61.      def __predict_tree_op(self,ip_arr,cols,mdf,get_prob = 'N'):
62.
63.          '''''Predicts lables for given tree for input array of attributes.
64.
65.          Parameters:
66.              Attributes Array,Meta Data frame
67.
68.          Summary:
69.              -Gets all Leaf Node conditions from Meta data frame
70.              -
   Checks which condition is met by test data and maps leaf node pertaining to the condition.
71.              -Predicts the most probable output class the probabilities in leaf node.
72.
73.          Returns:
74.              Predicted lables for the tree passed'''
75.
76.          op_lst = []
77.          leaf_node_df = mdf [mdf['leaf_f'] == 'Y']
78.          leaf_conditions = leaf_node_df['split_condition']
79.          for elm in ip_arr:
80.              p = pd.Series(elm,index = cols)
81.              for num,condition in enumerate(leaf_conditions):
82.                  check = []
83.                  for cond in condition:
84.                      col = cond.split()[0]
```

```
85.                          val = ' '.join(cond.split()[1:])
86.                          if  eval(str(p[col]) + val):
87.                              check.append('True')
88.                          else:
89.                              check.append('False')
90.                      if 'False' not in check:
91.                          class_prob_df = leaf_node_df.iloc[num,][self.data[self.lable_name].unique().to
      list()]
92.                          if get_prob == 'Y':
93.                              op = class_prob_df.values.tolist()
94.                              op_lst.append(op)
95.                          else:
96.                              op = class_prob_df[class_prob_df == class_prob_df.max()].index.values #ret
      urns column name with max value
97.                              op_lst.append(op[0])
98.                          break
99.
100.                 return np.array(op_lst)
101.
102.
103.             def predict(self,ip_arr):
104.
105.                 '''''Predicts lables from output lables obtained from each tree.
106.
107.                 Parameters:
108.                     Attributes Array,
109.
110.                 Returns:
111.                     Predicted lables'''
112.
113.                 votes_arr = []
114.
115.                 for tree in range(self.total_trees):
116.
117.                     cols = random_forest.feature_subseted[tree]
118.                     idx_cols = random_forest.features_idx_subseted[tree]
119.
120.                     idx_data =np.array([list(arr[idx_cols]) for arr in ip_arr])
121.
122.                     op = self.__predict_tree_op(idx_data,cols,random_forest.meta_dfs[tree])
123.                     votes_arr.append(op)
124.                     odf = pd.DataFrame(np.array(votes_arr))
125.                 return odf.mode().iloc[0].values
126.
127.             def write_meta_data(self,file_name):
128.
129.                 with pd.ExcelWriter(file_name) as workbook:
130.                     for tree,mdf in random_forest.meta_dfs.items():
131.                         mdf.to_excel(workbook, sheet_name='Tree_'+str(tree+1),index=False)
132.                 print("\nCheck '{}' for meta information at each tree.".format(file_name))
133.
134.
135.
136.
137.
138.             def predict_prob(self,ip_arr):
139.
140.                 '''''Predicts lables from output lables obtained from each tree.
141.
142.                 Parameters:
143.                     Attributes Array,
144.
145.                 Returns:
146.                     Predicted lables'''
147.
148.                 votes_arr = []
149.
150.                 for tree in range(self.total_trees):
151.
152.                     cols = random_forest.feature_subseted[tree]
153.                     idx_cols = random_forest.features_idx_subseted[tree]
154.
155.                     idx_data =np.array([list(arr[idx_cols]) for arr in ip_arr])
```

```
156.
157.                op = self.__predict_tree_op(idx_data,cols,random_forest.meta_dfs[tree],get_prob
     ='Y')
158.                votes_arr.append(op)
159.            prob = np.array(votes_arr)
160.            return np.mean(prob,axis=0)
```

## C4_5.PY

```
1.   import pandas as pd
2.   import numpy as np
3.
4.   from Algorithm.Utilities.Calculations import *
5.
6.
7.   class DT_classifier():
8.
9.       # Meta Data Frame to store decision tree information
10.      mdf  = pd.DataFrame()
11.
12.      def __init__(self,max_depth,train_df,feature_list,lable_name,lable_list):
13.          self.max_depth = max_depth
14.          self.train_df = train_df
15.          self.feature_list = feature_list
16.          self.lable_name = lable_name
17.          self.lable_list = lable_list
18.
19.
20.      def __tree_initial_split(self,feature_list,lable,df_test):
21.
22.          """Initiate initial tree split to get data into meta df for further spliting.
23.
24.          Parameters -
25.
26.              Attributes,Label name,Data
27.
28.          Returns -
29.
30.              Meta Data Frame
31.          """
32.          split_feature,threshold,overall_entropy = discretize(df_test,lable,feature_list)
33.          cols = ['node_number','parent_node','split_feature','split_category','split_condition','ov
     erall_entropy','leaf_f','next_split_features']
34.          val = [['1',np.nan,split_feature,np.nan,[threshold],overall_entropy,'N']]
35.          feature_for_next_split = feature_list.copy()
36.          feature_for_next_split.remove(split_feature)  # removing current split feature for further
      spliting
37.          val[0].append(feature_for_next_split)
38.          lables = df_test[lable].unique().tolist()
39.          for i in lables:
40.              cols.append(str(i))   # creating columns as per no of labels in op class
41.              class_data = len(df_test[df_test[lable] == i])
42.              total_data = len(df_test)
43.              val[0].append(class_data/total_data) # probability for that lable in the node
44.          meta_df = pd.DataFrame(val,columns=cols)
45.          #for num,levels in enumerate(df_test[split_feature].unique()): #Branches from the first no
     de
46.          for num,operator in enumerate([' <= ',' > ']):
47.              condition = operator +str(threshold)
48.              vals = [['1_'+str(num+1),'1',split_feature,condition,[split_feature+condition],np.nan,
     np.nan,feature_for_next_split]]
49.              for lable in lables: vals[0].append(np.nan) #appending values for op class probability

50.              meta_df = meta_df.append(pd.DataFrame(vals,columns=cols),ignore_index=True)[cols] # ad
     ding observations to meta_df
51.          return meta_df
52.
53.
```

```python
54.
55.
56.
57.     def __build_tree(self,data,meta_df,lable_name,lable_list,final_split = 'N'):
58.
59.         '''''Split tree from the provided state to next state.
60.
61.         Parameters:
62.
63.             Data,Meta Data Frame,Label name,Classes in lable,Final split flag
64.
65.         Summary:
66.
67.             -Filters the nodes that are not split using meta info
68.             -Filtered Nodes which are non-terminal will be split further.
69.             -
    Filtered Nodes are updated with Leaf flag,overall entropy and class probability information in met
    a data frame.
70.             -
    If final split flag = 'Y', only weights are updated no further spliting (If tree_depth has reached
     but node is not leaf then we do this as last step)
71.
72.         Returns:
73.             Temporary Data Frame containg meta information of new split done'''
74.
75.         df_node_to_split = meta_df[(meta_df['parent_node'].notnull()) & (pd.isnull(meta_df['overal
    l_entropy']))  ]
76.         cols = meta_df.columns.tolist()
77.         temp_df = pd.DataFrame(columns=cols)
78.         for i in range(0,len(df_node_to_split)):
79.             leaf = False
80.             series = df_node_to_split.iloc[i,]
81.             node_number,features_for_split,split_condition = series['node_number'],series['next_sp
    lit_features'],series['split_condition']
82.
83.             #Filtering data subset based on split condition dict
84.             subset = data
85.             for cond in split_condition:
86.                 key = cond.split()[0]
87.                 value = ' '.join(cond.split()[1:])
88.                 subset = eval('subset[subset[key]'+ value+']')
89.
90.             probs = []
91.             #Updating class label  propbailites and lead flag
92.             for lable in lable_list:
93.                 if len(subset) >0: # for zreo division error
94.                     prob = len(subset[subset[lable_name] == lable]) / len(subset)
95.                     meta_df.loc[meta_df['node_number'] ==series['node_number'],str(lable)] = prob

96.                     probs.append(prob)
97.             if (1 in probs) or (final_split == 'Y') or (len(features_for_split) == 0):
98.                 meta_df.loc[meta_df['node_number'] ==series['node_number'],'leaf_f'] = 'Y'
99.                 meta_df.loc[meta_df['node_number'] ==series['node_number'],'overall_entropy'] = 0

100.                 leaf = True
101.             else:
102.                 meta_df.loc[meta_df['node_number'] ==series['node_number'],'leaf_f'] = 'N'

103.
104.             if leaf or (final_split == 'Y') :
105.                 continue
106.
107.             #Get split feature for further split if its not a leaf node
108.
109.             split_feature,threshold,overall_entropy = discretize(subset,lable_name,features
    _for_split)
110.             meta_df.loc[meta_df['node_number'] ==series['node_number'],'overall_entropy'] =
     overall_entropy #Weighted entropy update
111.             next_features = features_for_split.copy()
112.             next_features.remove(split_feature) # child's next feature for split
113.
114.             #child observations entry
115.             for num,operator in enumerate([' <= ',' > ']):
```

```python
116.                   #for num,levels in enumerate(subset[split_feature].unique()):
117.                       condition = split_condition.copy()
118.                       condition.append(split_feature + operator + str(threshold))
119.                       vals = [[node_number+'_'+str(num+1),node_number,split_feature,operator + st
     r(threshold),condition,np.nan,np.nan,next_features]]
120.                       for lbl in lable_list: vals[0].append(np.nan)
121.                       temp_df = temp_df.append(pd.DataFrame(vals,columns=cols),ignore_index=True)
     [cols]
122.               return temp_df
123.
124.          def __recursive_tree(self,data,max_depth,lable_name,feature_list,lable_list):
125.
126.              '''''Initiates tree building and splits child nodes until convergence.
127.
128.              Parameters:
129.                  Data,Tree Depth,Label name,Data,Attributes,Output class names
130.
131.              Returns:
132.                  Meta Data Frame'''
133.
134.              meta_df = self.__tree_initial_split (feature_list,lable_name,data)
135.              max_depth = max_depth - 1 #as already 1 split has been done by us
136.              for i in range(0,max_depth):
137.                  temp_df = self.__build_tree(data,meta_df,lable_name,lable_list)
138.                  meta_df = meta_df.append(temp_df,ignore_index=True)[meta_df.columns.tolist()]
139.              temp_df = self.__build_tree(data,meta_df,lable_name,lable_list,final_split = 'Y') #
     just to update meta data no split occurs here
140.              meta_df = meta_df.append(temp_df,ignore_index=True)[meta_df.columns.tolist()]
141.              return meta_df
142.
143.
144.
145.
146.          def fit(self):
147.              df = self.__recursive_tree(self.train_df,self.max_depth,self.lable_name,self.featur
     e_list,self.lable_list)
148.              DT_classifier.mdf = df
```

---

## CALCULATIONS.PY

```python
1.   import numpy as np
2.
3.   def get_threshold(df,col,lable):
4.
5.       '''''Get best threshold value for the attribute
6.
7.           Parameters :
8.               Data frame,Attribute,Lable name
9.
10.          Returns:
11.              Entropy value,Threshold value'''
12.
13.      df = df[[col,lable]]
14.      uniq_vals = df[col].unique().tolist()
15.      uniq_vals.sort()
16.      thresholds = [(uniq_vals[idx]+uniq_vals[idx+1])/2 for idx in range(0,len(uniq_vals) -1)]
17.      weighted_entropy = []
18.
19.      for threshold in thresholds:
20.          left = df[df[col] <= threshold][lable]
21.          right = df[df[col] > threshold][lable]
22.          entropy = calc_wgtd_entropy_numeric(left,right)
23.          weighted_entropy.append(entropy)
24.      if thresholds == []:
25.          weighted_entropy = [0]
26.          thresholds = [0]
```

```python
27.
28.     return [min(weighted_entropy),thresholds[weighted_entropy.index(min(weighted_entropy))]]]
29.
30.
31.
32. def calc_wgtd_entropy_numeric( left, right):
33.
34.     '''''Get overall entropy for an attribute
35.
36.         Parameters :
37.             Data towards left and right of threshold (i.e, <= threshold & > threshold )
38.
39.         Returns:
40.             Overall entropy'''
41.
42.     total_elements = len(left) + len(right)
43.     ent_left = entropy(left)
44.     ent_right = entropy(right)
45.     weighted_entropy = ((len(left) / total_elements) * ent_left) + ((len(right) / total_elements)
   * ent_right)
46.     return weighted_entropy
47.
48.
49. def entropy(df):
50.
51.     '''''Get entropy for a subset
52.
53.         Parameters :
54.             Lable values of subset
55.
56.         Returns:
57.             Entropy'''
58.
59.     op_class, count = np.unique(df.values, return_counts=True)
60.     entropy = np.sum([(-
   count[i] / np.sum(count)) * np.log2(count[i] / np.sum(count)) for i in range(len(op_class))])
61.     return entropy
62.
63.
64.
65. def discretize(df,lable,feature_list):
66.
67.     '''''Discretization (binning) of continuous data
68.
69.         Parameters :
70.             Data Frame,Lable Name,Attributes
71.
72.         Returns:
73.             Attribute,Threshold,Overall Entropy'''
74.
75.     val = []
76.     d = {}
77.     for i in feature_list:
78.         op = get_threshold(df,i,lable)
79.         val.append(op)
80.         d[op[1]] = i
81.
82.     ent = [value[0] for value in val]
83.     thr = [value[1] for value in val]
84.     return d[thr[ent.index(min(ent))]],thr[ent.index(min(ent))],min(ent)
85.
86.
87. def get_bootstrap(data, bootstrap_size,random_seed):
88.
89.     ''''' Bootstraping for random forrest
90.
91.         Parameters :
92.             Data Frame, Bootstarp size, Random state
93.
94.         Returns:
95.             Boostrap sample'''
96.
97.     np.random.seed(random_seed)
```

```
98.     bootstraps = np.random.randint(low=0, high=len(data), size=bootstrap_size)
99.     df_bootstrap = data.iloc[bootstraps]
100.         return df_bootstrap
```

```
1.  import random
2.  import pandas as pd
3.  import numpy as np
4.  from Algorithm.Random_Forest import random_forest
5.  from sklearn.ensemble import RandomForestClassifier
6.  from sklearn.metrics import confusion_matrix
7.  from sklearn.model_selection import train_test_split
8.  from sklearn import metrics
9.  from sklearn.model_selection import cross_val_score,cross_val_predict
10.
11. df = pd.read_csv('Data\\hazzlenuts_preprocessed.csv')
12.
13. features = df.columns[0:-1].tolist()
14. output_lable_name = df.columns[-1]
15.
16.
17. x = df[features].values
18. y = df[output_lable_name].values
19.
20. x_train, x_test, y_train, y_test = train_test_split(x, y, test_size = 0.33, random_state = 1)
21.
22. train_df = pd.DataFrame(x_train,columns = features)
23. train_df[output_lable_name] = y_train
24. train_df
25.
26. def gen_random_splits(n):
27.     sk_rf_acc = []
28.     rf_acc = []
29.     for split in range(0,n):
30.
31.         print('\n---------------------------------Random split - {} -------------------
    '.format(split+1))
32.         x_train, x_test, y_train, y_test = train_test_split(x, y, test_size = 0.33, random_state =
    split)
33.
34.         train_df = pd.DataFrame(x_train,columns = features)
35.         train_df[output_lable_name] = y_train
36.         train_df
37.
38.         r = random_forest(data = train_df,total_trees = 11,tree_depth = 5,random_seed = 10,feature
    _per_tree = 6,bootstrap_size=100)
39.         r.fit()
40.         op= r.predict(x_test)
41.         acc_rf = metrics.accuracy_score(y_test, op)
42.         print('\nAccuracy for random_forest (implemented)= {}%'.format(round(acc_rf*100,2)))
43.         rf_acc.append(acc_rf)
44.
45.         classifier = RandomForestClassifier(n_estimators = 11,max_features = 6, criterion = 'entro
    py', random_state = 10,max_depth = 5)
46.         classifier.fit(x_train, y_train)
47.         op = classifier.predict(x_test)
48.         acc_sk = metrics.accuracy_score(y_test, op)
49.         print('\nAccuracy for random_forest (Sklearn)  = {}%'.format(round(acc_sk*100,2)))
50.         sk_rf_acc.append(acc_sk)
51.
52.     print('\n--------------------------------------------------------')
53.     print('\nMean Accuracy for random_forest (Implemented) after {} random splits  = {}%'.format(n
    ,round(np.array(rf_acc).mean()*100,3)))
```

```python
54.     print('\nMean Accuracy for random_forest (Sklearn) after {} random splits  = {}%'.format(n,rou
    nd(np.array(sk_rf_acc).mean()*100,3)))
55.
56. if __name__ == "__main__":
57.
58.     print('\n------------Random Forest (Implementation from scratch) for single random split ----
    --------------')
59.     r = random_forest(data = train_df,total_trees = 11,tree_depth = 5,random_seed = 10,feature_per
    _tree = 6,bootstrap_size=100)
60.
61.     r.fit()
62.     op= r.predict(x_test)
63.     prob_arr = r.predict_prob(x_test)
64.     pd.DataFrame({'Actual': y_test,'Predicted':op}).to_csv('Results\\Testing_Output.csv',index=Fal
    se)
65.
66.     print("\nCheck 'Results\\Testing_Output.csv' to compare actual and predicted values.")
67.     r.write_meta_data('Results\\Random_Forest_Meta_Data_Generated.xlsx')
68.
69.     print('\nConfusion Matrix \n')
70.     cm = confusion_matrix(y_test, op)
71.     print(pd.DataFrame(cm,columns=train_df.iloc[:,-1].unique(),index = train_df.iloc[:,-
    1].unique() ))
72.
73.
74.     acc = metrics.accuracy_score(y_test, op)
75.     print('\nAccuracy = {}%\n'.format(round(acc*100,2)))
76.
77.     print('-------------Random Forest (Implementation from Sklearn) for single random split ------
    --------------')
78.
79.
80.     classifier = RandomForestClassifier(n_estimators = 11,max_features = 6, criterion = 'entropy',
     random_state = 10,max_depth = 5)
81.     classifier.fit(x_train, y_train)
82.     op = classifier.predict(x_test)
83.
84.     print('\nConfusion Matrix\n')
85.     cm_sk = confusion_matrix(y_test, op)
86.     print(pd.DataFrame(cm_sk,columns=train_df.iloc[:,-1].unique(),index = train_df.iloc[:,-
    1].unique() ))
87.
88.     acc_sk = metrics.accuracy_score(y_test, op)
89.     print('\nAccuracy = {}%'.format(round(acc_sk*100,2)))
90.
91.
92.     print('\n\n<<<<<<<<<<<<<<<<<<<<<<<<<<<<<<<<<< Model Testing >>>>>>>>>>>>>>>>>>>>>>>>>>>>>>>>>>>
    >>\n')
93.     gen_random_splits(10)
94.
```

## CONSOLE OUTPUT:

Executing __main__.py file inside compressed zip.(Total execution time = 2m:45s [Inluding 10 random splits])

```
T:\NUIG DA\A2_Random_Forest>python __main__.py

-------------Random Forest (Implementation from scratch) for single random split -----------------

Check 'Results\Testing_Output.csv' to compare actual and predicted values.

Check 'Results\Random_Forest_Meta_Data_Generated.xlsx' for meta information at each tree.

Confusion Matrix

            c_cornuta  c_avellana  c_americana
c_cornuta        14           5            0
c_avellana        1          24            0
c_americana       2           0           21

Accuracy = 88.06%

-------------Random Forest (Implementation from Sklearn) for single random split -----------------

Confusion Matrix

            c_cornuta  c_avellana  c_americana
c_cornuta        19           0            0
c_avellana        1          23            1
c_americana       1           0           22

Accuracy = 95.52%

<<<<<<<<<<<<<<<<<<<<<<<<<<<<<<< Model Testing >>>>>>>>>>>>>>>>>>>>>>>>>>>>>>>

--------------------------------Random split - 1 ------------------

Accuracy for random_forest (implemented)= 86.57%

Accuracy for random_forest (Sklearn)  = 88.06%

--------------------------------Random split - 2 ------------------

Accuracy for random_forest (implemented)= 88.06%

Accuracy for random_forest (Sklearn)  = 95.52%

--------------------------------Random split - 3 ------------------

Accuracy for random_forest (implemented)= 86.57%

Accuracy for random_forest (Sklearn)  = 86.57%
```

```
-------------------------------Random split - 4 ------------------

Accuracy for random_forest (implemented)= 89.55%

Accuracy for random_forest (Sklearn)  = 88.06%

-------------------------------Random split - 5 ------------------

Accuracy for random_forest (implemented)= 91.04%

Accuracy for random_forest (Sklearn)  = 89.55%

-------------------------------Random split - 6 ------------------

Accuracy for random_forest (implemented)= 85.07%

Accuracy for random forest (Sklearn)  = 88.06%
-------------------------------Random split - 7 ------------------

Accuracy for random_forest (implemented)= 89.55%

Accuracy for random_forest (Sklearn)  = 89.55%

-------------------------------Random split - 8 ------------------

Accuracy for random_forest (implemented)= 89.55%

Accuracy for random_forest (Sklearn)  = 89.55%

-------------------------------Random split - 9 ------------------

Accuracy for random_forest (implemented)= 94.03%

Accuracy for random_forest (Sklearn)  = 95.52%

-------------------------------Random split - 10 ------------------

Accuracy for random_forest (implemented)= 86.57%

Accuracy for random_forest (Sklearn)  = 91.04%

----------------------------------------------------

Mean Accuracy for random_forest (Implemented) after 10 random splits  = 88.657%

Mean Accuracy for random_forest (Sklearn) after 10 random splits  = 90.149%
```