

Program – 12

AIM: To Implement Decision Tree classifier in python

Introduction and Theory

Decision tree induction is the learning of decision trees from class-labeled training tuples. A decision tree is a flowchart-like tree structure, where each internal node (nonleaf node) denotes a test on an attribute, each branch represents an outcome of the test, and each leaf node (or terminal node) holds a class label. The topmost node in a tree is the root node.

Given a tuple, X , for which the associated class label is unknown, the attribute values of the tuple are tested against the decision tree. A path is traced from the root to a leaf node, which holds the class prediction for that tuple. Decision trees can easily be converted to classification rules.

ID3, C4.5, and CART adopt a greedy (i.e., nonbacktracking) approach in which decision trees are constructed in a top-down recursive divide-and-conquer manner. Most algorithms for decision tree induction also follow a top-down approach, which starts with a training set of tuples and their associated class labels. The training set is recursively partitioned into smaller subsets as the tree is being built.

Algorithm: Generate_decision_tree. Generate a decision tree from the training tuples of data partition, D .

Input:

- Data partition, D , which is a set of training tuples and their associated class labels;
- *attribute_list*, the set of candidate attributes;
- *Attribute_selection_method*, a procedure to determine the splitting criterion that “best” partitions the data tuples into individual classes. This criterion consists of a *splitting_attribute* and, possibly, either a *split_point* or *splitting_subset*.

Output: A decision tree.

Method:

- (1) create a node N ;
- (2) **if** tuples in D are all of the same class, C , **then**
- (3) return N as a leaf node labeled with the class C ;
- (4) **if** *attribute_list* is empty **then**
- (5) return N as a leaf node labeled with the majority class in D ; // majority voting
- (6) apply *Attribute_selection_method*(D , *attribute_list*) to find the “best” *splitting_criterion*;
- (7) label node N with *splitting_criterion*;
- (8) **if** *splitting_attribute* is discrete-valued **and**
 multiway splits allowed **then** // not restricted to binary trees
- (9) *attribute_list* \leftarrow *attribute_list* – *splitting_attribute*; // remove *splitting_attribute*
- (10) **for each** outcome j of *splitting_criterion*
 // partition the tuples and grow subtrees for each partition
- (11) let D_j be the set of data tuples in D satisfying outcome j ; // a partition
- (12) **if** D_j is empty **then**
- (13) attach a leaf labeled with the majority class in D to node N ;
- (14) **else** attach the node returned by *Generate_decision_tree*(D_j , *attribute_list*) to node N ;
- endfor**
- (15) return N ;

An attribute selection measure is a heuristic for selecting the splitting criterion that “best” separates a given data partition, D , of class-labeled training tuples into individual classes. If we were to split D into smaller partitions according to the outcomes of the splitting criterion, ideally each partition would be pure (i.e., all the tuples that fall into a given partition would belong to the same class). Conceptually, the “best” splitting criterion is the one that most

Program – 12

closely results in such a scenario. Attribute selection measures are also known as splitting rules because they determine how the tuples at a given node are to be split.

ID3 uses Information gain:

$$\begin{aligned} Info(D) &= - \sum_{i=1}^m p_i \log_2(p_i) \\ Info_A(D) &= \sum_{j=1}^v \frac{|D_j|}{|D|} \times Info(D_j) \\ Gain(A) &= Info(D) - info_A(D) \end{aligned}$$

C4.5 uses gain ratio

$$\begin{aligned} SplitInfo_A(D) &= - \sum_{j=1}^v \frac{|D_j|}{|D|} \times \log_2 \frac{|D_j|}{|D|} \\ GainRatio(A) &= \frac{Gain(A)}{SplitInfo_A(D)} \end{aligned}$$

CART uses the GINI index

$$\begin{aligned} Gini(D) &= 1 - \sum_{i=1}^m p_i^2 \\ Gini_A(D) &= \frac{|D_1|}{|D|} Gini(D_1) + \frac{|D_2|}{|D|} Gini(D_2) \\ \Delta Gini(A) &= Gini(D) - Gini_A(D) \end{aligned}$$

Strengths and Weakness of Decision Tree approach

Strengths:

- Decision trees are able to generate understandable rules.
- Decision trees perform classification without requiring much computation.
- Decision trees are able to handle both continuous and categorical variables.
- Decision trees provide a clear indication of which fields are most important for prediction or classification.

Weaknesses:

- Decision trees are less appropriate for estimation tasks where the goal is to predict the value of a continuous attribute.
- Decision trees are prone to errors in classification problems with many class and relatively small number of training examples.
- Decision tree can be computationally expensive to train. The process of growing a decision tree is computationally expensive. At each node, each candidate splitting field must be sorted before its best split can be found. In some algorithms, combinations of fields are used and a search must be made for optimal combining weights. Pruning algorithms can also be expensive since many candidate sub-trees must be formed and compared.

Program – 12

Code

```
1  # importing libraries
2  import pandas as pd
3  import numpy as np
4  from pprint import pprint
5
6  # calculates the entropy
7  def entropy(target_col):
8      elements, counts = np.unique(target_col, return_counts=True)
9      result = np.sum([(
10 counts[i]/np.sum(counts))*np.log2(counts[i]/np.sum(counts)) for i in
11 range(len(elements))])
12     return result
13
14 # calculate the information gain = Gain(A) / SplitInfo_A(D)
15 def information_gain(data, split_attribute_name,
16 target_name="class"):
17     """
18     calculates the information gain of the dataset for attribute A
19     :param data: the dataset for which we're calculating the
20 information gain
21     :param split_attribute_name: name of the splitting attribute.
22     :param target_name: name of the target / class.
23     :return: the information gain
24     """
25     # Calculate the entropy of the total dataset
26     total_entropy = entropy(data[target_name])
27
28     # Calculate the values and the corresponding counts for the
29 split attribute
30     vals, counts = np.unique(data[split_attribute_name],
31 return_counts=True)
32     # Calculate the weighted entropy
33     weighted_entropy = np.sum(
34         [(counts[i] / np.sum(counts)) *
35 entropy(data.where(data[split_attribute_name] ==
36 vals[i]).dropna()[target_name])
37         for i in range(len(vals))])
38     # Calculate the information gain
39     information_gain_ = total_entropy - weighted_entropy
40     return information_gain_
41
42 # the ID3 Decision tree algorithm
43 def id3(data, original_data, features,
44 target_attribute_name="class", parent_node_class=None):
45     """
46     id3 algorithm
47     :param data: the dataset for which we're running the id3
48 algorithm
49     :param original_data: the original dataset needed to calculate
50 the mode target feature value of the original dataset
51     in the case the dataset delivered by the first parameter is
52 empty
53     :param features: the feature space of the dataset . This is
54 needed for the recursive call since during the tree
55
```

Program – 12

```
56     growing process we have to remove features from our dataset ->
57 Splitting at each node
58     :param target_attribute_name: the name of the target attribute
59     :param parent_node_class: This is the value or class of the mode
60 target feature value of the parent node for a
61     specific node. This is also needed for the recursive call since
62 if the splitting leads to a situation that there are
63     no more features left in the feature space, we want to return
64 the mode target feature value of the direct parent node.
65
66     :return: The learnt decision tree
67     """
68     # Define the stopping criteria -> If one of this is satisfied,
69 we want to return a leaf node
70     # If all target_values have the same value, return this value
71     if len(np.unique(data[target_attribute_name])) <= 1:
72         return np.unique(data[target_attribute_name])[0]
73     # If the dataset is empty, return the mode target feature value
74 in the original dataset
75     elif len(data) == 0:
76         return np.unique(original_data[target_attribute_name])[
77
78 np.argmax(np.unique(original_data[target_attribute_name],
79 return_counts=True)[1])]
80     # If the feature space is empty, return the mode target feature
81 value of the direct parent node -> Note that
82     # the direct parent node is that node which has called the
83 current run of the ID3 algorithm and hence
84     # the mode target feature value is stored in the
85 parent_node_class variable.
86     elif len(features) == 0:
87         return parent_node_class
88     # If none of the above holds true, grow the tree!
89     else:
90         # Set the default value for this node -> The mode target
91 feature value of the current node
92         parent_node_class = np.unique(data[target_attribute_name])[
93         np.argmax(np.unique(data[target_attribute_name],
94 return_counts=True)[1])]
95         # Select the feature which best splits the dataset
96         item_values = [information_gain(data, feature,
97 target_attribute_name) for feature in
98         features] # Return the information gain
99 values for the features in the dataset
100         best_feature_index = np.argmax(item_values)
101         best_feature = features[best_feature_index]
102         # Create the tree structure. The root gets the name of the
103 feature (best_feature) with the maximum information
104         # gain in the first run
105         tree = {best_feature: {}}
106         # Remove the feature with the best information gain from the
107 feature space
108         features = [i for i in features if i != best_feature]
109         # Grow a branch under the root node for each possible value
110 of the root node feature
111         for value in np.unique(data[best_feature]):
112             value = value
```

Program – 12

```
113         # Split the dataset along the value of the feature with
114 the largest
115         # information gain and then create sub_datasets
116         sub_data = data.where(data[best_feature] ==
117 value).dropna()
118
119         # Call the ID3 algorithm for each of those sub_datasets
120 with the
121         # new parameters -> Here the recursion comes in!
122         subtree = id3(sub_data, dataset, features,
123 target_attribute_name, parent_node_class)
124
125         # Add the sub tree, grown from the sub_dataset to the
126 tree under the root node
127         tree[best_feature][value] = subtree
128
129     return tree
130
131
132 def predict(query, tree, default=1):
133     """
134     prediction function for new unseen data
135     :param query: a dictionary entry, {"feature_name" : value, ... ,
136 "feature_name" : value}
137     :param tree: the ID3 tree
138     :param default: base value
139     :return: predicted class
140     """
141     # 1. Check for every feature in the query instance if this feature is
142 existing in the tree.keys() for the first call, tree.keys() only contains
143 the value for the root node -> if this value is not existing, we can not
144 make a prediction and have to return the default value which is the majority
145 value of the target feature
146     for key in list(query.keys()):
147         if key in list(tree.keys()):
148             # 2. First of all we have to take care of a important
149 fact: Since we train our model with a database A and then show our
150 model a unseen query it may happen that the feature values of these
151 query are not existing in our tree model because non of the training
152 instances has had such a value for this specific feature.
153             # For instance imagine the situation where your model
154 has only seen animals with one to four legs - The "legs" node in
155 your model will only have four outgoing branches (from one to four).
156 If you now show your model a new instance (animal) which has for the
157 legs feature the vale 5, you have to tell your model what to do in
158 such a situation because otherwise there is no classification
159 possible
160             # because in the classification step you try to run
161 down the outgoing branch with the value 5 but there is no such a
162 branch. Hence: Error and no Classification! We can address this
163 issue with a classification value of for instance (999) which tells
164 us that there is no classification possible or we assign the most
165 frequent target feature value of our dataset used to train the
166 model. Or, in for instance medical application we can return the
167 most worse case - just to make sure... We can also return the most
168 frequent value of the direct parent node. To make a long story
169 short, we have to tell the model what to do in this situation. In
```

Program – 12

```
170 our example, since we are dealing with animal species where a false
171 classification is not that critical, we will assign the value 1
172 which is the value for the mammal species (for convenience).
173     try:
174         result = tree[key][query[key]]
175     except:
176         return default
177
178     # 3. Address the key in the tree which fits the value
179     for key --> Note that key == the features in the query. Because we
180     want the tree to predict the value which is hidden under the key
181     value (imagine you have a drawn tree model on the table in front of
182     you and you have a query instance for which you want to predict the
183     target feature - What are you doing? - Correct: You start at the
184     root node and wander down the tree comparing your query to the node
185     values. Hence you want to have the value which is hidden under the
186     current node. If this is a leaf, perfect, otherwise you wander the
187     tree deeper until you get to a leaf node. Though, you want to have
188     this "something" [either leaf or sub_tree] which is hidden under the
189     current node and hence we must address the node in the tree which ==
190     the key value from our query instance. This is done with
191     tree[keys]. Next you want to run down the branch of this node which
192     is equal to the value given "behind" the key value of your query
193     instance e.g. if you find "legs" == to tree.keys() that is, for the
194     first run == the root node. You want to run deeper and therefore you
195     have to address the branch at your node whose value is == to the
196     value behind key. This is done with query[key]
197     #     e.g. query[key] == query['legs'] == 0 --> Therewith
198     we run down the branch of the node with the value 0. Summarized, in
199     this step we want to address the node which is hidden behind a
200     specific branch of the root node (in the first run) this is done
201     with: result = [key][query[key]]
202     result = tree[key][query[key]]
203     # 4. As said in the 2. step, we run down the tree along
204     nodes and branches until we get to a leaf node. That is, if result
205     = tree[key][query[key]] returns another tree object (we have
206     represented this by a dict object -> that is if result is a dict
207     object) we know that we have not arrived at a root node and have to
208     run deeper the tree. Okay... Look at your drawn tree in front of
209     you... what are you doing?...well, you run down the next branch...
210     exactly as we have done it above with the slight difference that we
211     already have passed a node and therewith have to run only a fraction
212     of the tree --> You clever guy! That "fraction of the tree" is
213     exactly what we have stored under 'result'. So we simply call our
214     predict method using the same query instance (we do not have to drop
215     any features from the query instance since for instance the feature
216     for the root node will not be available in any of the deeper
217     sub_trees and hence we will simply not find that feature) as well as
218     the "reduced / sub_tree" stored in result.
219     if isinstance(result, dict):
220         return predict(query, result)
221     else:
222         return result
223     # splitting the data into test-train splits for checking performance
224     on unseen data
225     def train_test_split(dataset):
226
```

Program – 12

```
227     training_data = dataset.iloc[:80].reset_index(drop=True) # We
228 drop the index respectively relabel the index
229     # starting form 0, because we do not want to run into errors
230 regarding the row labels / indexes
231     testing_data = dataset.iloc[80:].reset_index(drop=True)
232     return training_data, testing_data
233
234 # testing the tree model, get prediction accuracy
235 def test(data, tree):
236     # Create new query instances by simply removing the target
237 feature column from the original dataset and
238     # convert it to a dictionary
239     queries = data.iloc[:, :-1].to_dict(orient="records")
240
241     # Create a empty DataFrame in whose columns the prediction of
242 the tree are stored
243     predicted = pd.DataFrame(columns=["predicted"])
244
245     # Calculate the prediction accuracy
246     for i in range(len(data)):
247         predicted.loc[i, "predicted"] = predict(queries[i], tree,
248 1.0)
249     print('The prediction accuracy is: ',
250 (np.sum(predicted["predicted"] == data["class"]) / len(data)) * 100,
251 '%')
252
253
254
255 """
256 Train the tree, Print the tree and predict the accuracy
257 """
258
259 if __name__ == '__main__':
260     # loading the dataset
261     dataset = pd.read_csv('zoo.csv', names=['animal_name', 'hair',
262 'feathers', 'eggs', 'milk',
263                                     'airbone', 'aquatic',
264 'predator', 'toothed', 'backbone',
265                                     'breathes', 'venomous',
266 'fins', 'legs', 'tail', 'domestic', 'catsize',
267                                     'class'])
268     print(dataset.head(10))
269     # dropping the class column
270     dataset = dataset.drop('animal_name', axis=1)
271     print(dataset.head(10))
272     # splitting data
273     training_data = train_test_split(dataset)[0]
274     testing_data = train_test_split(dataset)[1]
275     # training the tree
276     tree = id3(training_data, training_data,
277 training_data.columns[:-1])
278     # printing the learnt tree in the form of a dictionary
279     pprint(tree)
280     # get test performance of the tree
281     test(testing_data, tree)
```

Program – 12

Results and Outputs:

```
DWDM_LAB — -bash — 128x41

[(ML) Anurags-MacBook-Air:DWDM_LAB jarvis$ python DecisionTree.py
animal_name hair feathers eggs milk airborne aquatic ... venomous fins legs tail domestic catsize class
0 aardvark 1 0 0 1 0 0 ... 0 0 4 0 0 1 1
1 antelope 1 0 0 1 0 0 ... 0 0 4 1 0 1 1
2 bass 0 0 1 0 0 1 ... 0 1 0 1 0 0 4
3 bear 1 0 0 1 0 0 ... 0 0 4 0 0 1 1
4 boar 1 0 0 1 0 0 ... 0 0 4 1 0 1 1
5 buffalo 1 0 0 1 0 0 ... 0 0 4 1 0 1 1
6 calf 1 0 0 1 0 0 ... 0 0 4 1 1 1 1
7 carp 0 0 1 0 0 1 ... 0 1 0 1 1 0 4
8 catfish 0 0 1 0 0 1 ... 0 1 0 1 0 0 4
9 cavy 1 0 0 1 0 0 ... 0 0 4 0 1 0 1

[10 rows x 18 columns]
hair feathers eggs milk airborne aquatic predator ... venomous fins legs tail domestic catsize class
0 1 0 0 1 0 0 1 ... 0 0 4 0 0 1 1
1 1 0 0 1 0 0 0 ... 0 0 4 1 0 1 1
2 0 0 1 0 0 1 1 ... 0 1 0 1 0 0 4
3 1 0 0 1 0 0 1 ... 0 0 4 0 0 1 1
4 1 0 0 1 0 0 1 ... 0 0 4 1 0 1 1
5 1 0 0 1 0 0 0 ... 0 0 4 1 0 1 1
6 1 0 0 1 0 0 0 ... 0 0 4 1 1 1 1
7 0 0 1 0 0 1 0 ... 0 1 0 1 1 0 4
8 0 0 1 0 0 1 1 ... 0 1 0 1 0 0 4
9 1 0 0 1 0 0 0 ... 0 0 4 0 1 0 1

[10 rows x 17 columns]
{'legs': {0: {'fins': {0.0: {'toothed': {0.0: 7.0, 1.0: 3.0}},
1.0: {'eggs': {0.0: 1.0, 1.0: 4.0}}}},
2: {'hair': {0.0: 2.0, 1.0: 1.0}},
4: {'hair': {0.0: {'toothed': {0.0: 7.0, 1.0: 5.0}}, 1.0: 1.0}},
6: {'aquatic': {0.0: 6.0, 1.0: 7.0}},
8: 7.0}}
The prediction accuracy is: 85.71428571428571 %
[(ML) Anurags-MacBook-Air:DWDM_LAB jarvis$
[(ML) Anurags-MacBook-Air:DWDM_LAB jarvis$
[(ML) Anurags-MacBook-Air:DWDM_LAB jarvis$
[(ML) Anurags-MacBook-Air:DWDM_LAB jarvis$
```

Findings and Learnings:

1. We have Implemented Decision tree through the ID3 algorithm in python 3.
2. We have learned the nuances of the Decision tree learning.
3. We have learnt about the applications, strengths and weaknesses of Decision tree Learning.