

Design and Implementation of Iterative Dichotomiser 3 (ID3) Decision Tree Classifier with Python

Sami Alperen Akgun
sami.alperen.akgun@gmail.com

I. ID3 BASED ON INFORMATION GAIN

A. Implementation

Iterative Dichotomiser 3 (ID3) decision tree classifier based on maximizing information gain (IG) was implemented for this question in Python programming language using numpy library. The implemented functions *BuildingDT* and *BuildingDT_{cont}* for data with discrete attributes and data with continuous attributes respectively can be seen below:

```
1 def BuildingDT(input_X, input_Y, feature_X_values, feature_Y_values, Attribute, mc_label):
2     """
3     This function creates an ID3 decision tree
4     Inputs: input_X --> np array contains data with features/attributes
5             input_Y --> np array contains output labels
6             feature_X_values --> list composed of possible values of features
7             feature_Y_values --> list composed of possible values of out labels
8             Attribute --> selected attribute/feature number(s) (column index)
9                     It is a list composed of column indices
10            mc_label --> most common label in the original data Y
11     Output: Decision Tree DT (class)
12     """
13     # If there is no output label
14     if input_Y.shape[0] == 0:
15         return DT(LeafValue=mc_label)
16
17     children = [] #children nodes from the current selected node
18     mcv = most_common_value(input_Y, feature_Y_values) # most common value of given input Y
19
20     # If all labels of input_Y is the same or Attribute is empty
21     all_values_same_boolean = check_all_values_same(input_Y, feature_Y_values)
22     if all_values_same_boolean or Attribute == []:
23         return DT(LeafValue=mcv)
24     else:
25         best_IG = 0
26         best_set_X = []
27         best_set_Y = []
28         best_Attribute = 0
29         current_entropy = entropy(input_Y, feature_Y_values)
30         for index in Attribute:
31             # index is the column index --> decides which attribute to use
32             divided_X, divided_Y = divide_data(input_X, input_Y, index, feature_X_values)
33
34             for j in range(len(divided_Y)): #iterate through each possible feature value
35                 # For tic-tac-toe --> divided_Y contains subsets for "x", "o" and "b"
36                 current_IG = information_gain(divided_Y[j], feature_Y_values, current_entropy)
37                 if current_IG >= best_IG:
38                     best_IG = current_IG
39                     best_Attribute = index
40                     best_feature_ind = j
41
42         # Now best attribute is selected --> divide data using that attribute
43         # best_Attribute --> gives the column number of best attribute
44         # best_feature_ind --> gives which feature value has the most IG, i.e. "x", "o" or "b" in tic-tac
45         best_set_X, best_set_Y = divide_data(input_X, input_Y, best_Attribute, feature_X_values)
46
47         # Remove the current attribute from the attribute list to
48         # create remaning attributes --> no need to use same attribute over and
49         # over again! --> If it is needed, it is gonna show up in the recursion anyway
50         remaning_Attributes = []
51         for index in Attribute:
52             if index != best_Attribute:
53                 remaning_Attributes.append(index)
```

```

53         if index != best_Attribute:
54             remaning_Attributes.append(index)
55
56
57     for i in range(len(best_set_X)):
58         children.append( BuildingDT(best_set_X[i],best_set_Y[i],feature_X_values,feature_Y_values,
59                                 remaning_Attributes,mcv) )
60
61     return DT (FeatureNum=best_Attribute,Children=children)

1 def BuildingDT_cont (input_X,input_Y,feature_X_values,feature_Y_values,Attribute,mc_label):
2     """
3     This function creates an ID3 decision tree
4     Inputs: input_X --> np array contains data with features/attributes
5             input_Y --> np array contains output labels
6             feature_X_values --> list composed of possible values of features
7             feature_Y_values --> list composed of possible values of out labels
8             Attribute --> selected attribute/feature number(s) (column index)
9                     It is a list composed of column indices
10            mc_label --> most common label in the original data Y
11    Output: Decision Tree DT (class)
12    """
13    # If there is no output label
14    if input_Y.shape[0] == 0:
15        return DT_cont (LeafValue=mc_label)
16
17    left = [] # contains values less or equal than threshold
18    right = [] # higher than threshold
19    mcv = most_common_value(input_Y,feature_Y_values) # most common value of given input Y
20
21    # If all labels of input_Y is the same or Attribute is empty
22    all_values_same_boolean = check_all_values_same(input_Y,feature_Y_values)
23    if all_values_same_boolean or Attribute == []:
24        return DT_cont (LeafValue=mcv)
25    else:
26        best_IG = 0
27        best_set_X = []
28        best_set_Y = []
29        best_Attribute = 0
30        best_threshold = 0.0
31        current_entropy = entropy(input_Y,feature_Y_values)
32        for index in Attribute:
33            # index is the column index --> decides which attribute to use
34            all_thresholds = return_all_thresholds(input_X,input_Y,index)
35            for threshold in all_thresholds:
36                # Get divided data for each threshold for a given attribute = index
37                divided_X, divided_Y = divide_data(input_X,input_Y,index,feature_X_values,threshold)
38
39                for j in range(len(divided_Y)): #iterate through each possible feature value
40                    # For wine.data --> divided_Y contains subsets for lower then threshold
41                    # and higher than the threshold
42                    current_IG = information_gain(divided_Y[j],feature_Y_values,current_entropy)
43                    if current_IG >= best_IG:
44                        best_IG = current_IG
45                        best_Attribute = index
46                        best_feature_ind = j
47                        best_threshold = threshold
48
49
50    # Now best attribute and threshold are selected --> divide data using that attribute and thresh
51    # best_Attribute --> gives the column number of best attribute
52    # best threshold --> gives the best threshold value which gives max IG
53    # best_feature_ind --> gives which feature value has the most IG, i.e. "x", "o" or "b" in tic-tac
54    best_set_X, best_set_Y = divide_data(input_X,input_Y,best_Attribute,feature_X_values,best_threshold
55    )
56
57    # Remove the current attribute from the attribute list to
58    # create remaning attributes --> no need to use same attribute over and
59    # over again! --> If it is needed, it is gonna show up in the recursion anyway
60    remaning_Attributes = []
61    for index in Attribute:
62        if index != best_Attribute:
63            remaning_Attributes.append(index)

```

```

64     # Left contains values less or equal than threshold
65     # Right contains values more than threshold
66     # For more detail, check divide data function
67
68     return DT_cont (Left=BuildingDT_cont (best_set_X[0],best_set_Y[0],feature_X_values,feature_Y_values,
69                                     remaning_Attributes,mcv),
70                    Right=BuildingDT_cont (best_set_X[1],best_set_Y[1],feature_X_values,feature_Y_values,
71                                     remaning_Attributes,mcv),
72                    FeatureNum=best_Attribute,ThresholdNum=best_threshold
73                    )

```

B. Training and Performance

10-times-10-fold cross validation technique was employed to divide the data into training and test set. Decision trees were created using training data and accuracy of them were calculated using test data.

Here are the obtained results for tic-tac-toe endgame dataset:

```

1 Mean of the accuracy:  0.78875 --> 78%
2 Var of the accuracy:  2.7769225320698284e-07
3 Best accuracy:  0.8854166666666666 --> 89%

```

Here are the obtained results for wine dataset:

```

1 Mean of the accuracy:  0.7816666666666667 --> 78%
2 Var of the accuracy:  1.4005687014174652e-05
3 Best accuracy:  0.9444444444444444 --> 94%

```

Finally, confusion matrices were calculated for the tree which gives the best accuracy among 10-times-10-fold cross validation. It can be seen in Table I, where rows shows actual classes and columns show predicted classes.

TABLE I
CONFUSION MATRICES

Predicted Numbers	
Positive	Negative
60	5
6	25

(a) Confusion Matrix
for tic-tac-toe endgame
Dataset

Predicted Numbers		
Class "1"	Class "2"	Class "3"
6	0	0
0	7	0
0	1	4

(b) Confusion Matrix for wine Dataset

Average mean of accuries for both dataset are almost the same, but variance in wine dataset is lower than tic-tac-toe dataset, which shows that decision tree trained for wine dataset is more successful than the other one since result of 10-times-10-fold cross validation don't change too much. If we look at confusion matrices, we can infer that classifier for tic-tac-toe dataset has balanced errors, i.e. the wrong prediction numbers are almost the same for positive and negative label case. On the other hand, confusion matrix for wine dataset clearly shows that trained decision tree is almost perfect and it predicted only one example wrong, which was normally belong to class 3, but tree predicted as class 2. One last thing to note that, since number of test samples are quite low for wine dataset, even having one wrong prediction reduces the accuracy a lot in a percentage wise.

II. ID3 BASED ON GAIN RATIO

A. Implementation

Same procedure was executed once again, but this time decision trees were created based on gain ratio (GR) instead of information gain (IG). One can see implemented functions *BuildingDT_{GR}* and *BuildingDT_{cont}GR* for data with discrete attributes and data with continious attributes below:

```

1 def BuildingDT_GR(input_X,input_Y,feature_X_values,feature_Y_values,Attribute,mc_label):
2     """
3     This function creates an ID3 decision tree based on maximizing gain ratio (GR)
4     Inputs: input_X --> np array contains data with features/attributes
5             input_Y --> np array contains output labels
6             feature_X_values --> list composed of possible values of features
7             feature_Y_values --> list composed of possible values of out labels
8             Attribute --> selected attribute/feature number(s) (column index)
9                     It is a list composed of column indices
10            mc_label --> most common label in the original data Y
11     Output: Decision Tree DT (class)
12     """

```

```

13 # If there is no output label
14 if input_Y.shape[0] == 0:
15     return DT(LeafValue=mc_label)
16
17 children = [] #children nodes from the current selected node
18 mcv = most_common_value(input_Y,feature_Y_values) # most common value of given input Y
19
20 # If all labels of input_Y is the same or Attribute is empty
21 all_values_same_boolean = check_all_values_same(input_Y,feature_Y_values)
22 if all_values_same_boolean or Attribute == []:
23     return DT(LeafValue=mcv)
24 else:
25     best_GR = 0 #best gain ratio
26     best_set_X = []
27     best_set_Y = []
28     best_Attribute = 0
29     current_entropy = entropy(input_Y,feature_Y_values)
30     for index in Attribute:
31         # index is the column index --> decides which attribute to use
32         divided_X, divided_Y = divide_data(input_X,input_Y,index,feature_X_values)
33
34         for j in range(len(divided_Y)): #iterate through each possible feature value
35             # For tic-tac-toe --> divided_Y contains subsets for "x", "o" and "b"
36             current_IG = information_gain(divided_Y[j],feature_Y_values,current_entropy)
37             current_GR = gain_ratio(divided_Y[j],feature_Y_values,current_IG)
38             if current_GR >= best_GR:
39                 best_GR = current_GR
40                 best_Attribute = index
41                 best_feature_ind = j
42
43
44     # Now best attribute is selected --> divide data using that attribute
45     # best_Attribute --> gives the column number of best attribute
46     # best_feature_ind --> gives which feature value has the most IG, i.e. "x", "o" or "b" in tic-tac
47     best_set_X, best_set_Y = divide_data(input_X,input_Y,best_Attribute,feature_X_values)
48
49     # Remove the current attribute from the attribute list to
50     # create remaning attributes --> no need to use same attribute over and
51     # over again! --> If it is needed, it is gonna show up in the recursion anyway
52     remaning_Attributes = []
53     for index in Attribute:
54         if index != best_Attribute:
55             remaning_Attributes.append(index)
56
57
58     for i in range(len(best_set_X)):
59         children.append( BuildingDT(best_set_X[i],best_set_Y[i],feature_X_values,feature_Y_values,
60                                     remaning_Attributes,mcv) )
61
62     return DT(FeatureNum=best_Attribute,Children=children)

```

```

1 def BuildingDT_cont_GR(input_X,input_Y, feature_X_values,feature_Y_values,Attribute,mc_label):
2     """
3     This function creates an ID3 decision tree based on maximizing gain ratio (GR)
4     Inputs: input_X --> np array contains data with features/attributes
5             input_Y --> np array contains output labels
6             feature_X_values --> list composed of possible values of features
7             feature_Y_values --> list composed of possible values of out labels
8             Attribute --> selected attribute/feature number(s) (column index)
9                     It is a list composed of column indices
10            mc_label --> most common label in the original data Y
11     Output: Decision Tree DT (class)
12     """
13     # If there is no output label
14     if input_Y.shape[0] == 0:
15         return DT_cont(LeafValue=mc_label)
16
17     left = [] # contains values less or equal than threshold
18     right = [] # higher than threshold
19     mcv = most_common_value(input_Y,feature_Y_values) # most common value of given input Y
20
21     # If all labels of input_Y is the same or Attribute is empty
22     all_values_same_boolean = check_all_values_same(input_Y,feature_Y_values)
23     if all_values_same_boolean or Attribute == []:

```

```

24     return DT_cont(LeafValue=mcv)
25 else:
26     best_GR = 0
27     best_set_X = []
28     best_set_Y = []
29     best_Attribute = 0
30     best_threshold = 0.0
31     current_entropy = entropy(input_Y, feature_Y_values)
32     for index in Attribute:
33         # index is the column index --> decides which attribute to use
34         all_thresholds = return_all_thresholds(input_X, input_Y, index)
35         for threshold in all_thresholds:
36             # Get divided data for each threshold for a given attribute = index
37             divided_X, divided_Y = divide_data(input_X, input_Y, index, feature_X_values, threshold)
38
39             for j in range(len(divided_Y)): #iterate through each possible feature value
40                 # For wine.data --> divided_Y contains subsets for lower then threshold
41                 # and higher than the threshold
42                 current_IG = information_gain(divided_Y[j], feature_Y_values, current_entropy)
43                 current_GR = gain_ratio(divided_Y[j], feature_Y_values, current_IG)
44                 if current_GR >= best_GR:
45                     best_GR = current_GR
46                     best_Attribute = index
47                     best_feature_ind = j
48                     best_threshold = threshold
49
50
51     # Now best attribute and threshold are selected --> divide data using that attribute and thresh
52     # best_Attribute --> gives the column number of best attribute
53     # best_threshold --> gives the best threshold value which gives max IG
54     # best_feature_ind --> gives which feature value has the most IG, i.e. "x", "o" or "b" in tic-tac
55     best_set_X, best_set_Y = divide_data(input_X, input_Y, best_Attribute, feature_X_values, best_threshold
56 )
57
58     # Remove the current attribute from the attribute list to
59     # create remaning attributes --> no need to use same attribute over and
60     # over again! --> If it is needed, it is gonna show up in the recursion anyway
61     remaning_Attributes = []
62     for index in Attribute:
63         if index != best_Attribute:
64             remaning_Attributes.append(index)
65
66     # Left contains values less or equal than threshold
67     # Right contains values more than threshold
68     # For more detail, check divide data function
69     #left.append( BuildingDT_cont(best_set_X[0], best_set_Y[0], feature_X_values, feature_Y_values,
70                                #remaning_Attributes, mcv) )
71     #right.append( BuildingDT_cont(best_set_X[1], best_set_Y[1], feature_X_values, feature_Y_values,
72                                #remaning_Attributes, mcv) )
73     return DT_cont(Left=BuildingDT_cont(best_set_X[0], best_set_Y[0], feature_X_values, feature_Y_values,
74                                       remaning_Attributes, mcv),
75                   Right=BuildingDT_cont(best_set_X[1], best_set_Y[1], feature_X_values, feature_Y_values,
76                                       remaning_Attributes, mcv),
77                   FeatureNum=best_Attribute, ThresholdNum=best_threshold
78 )

```

B. Training and Performance

Again 10-times-10-fold cross validation was executed to evaluate the implemented decision trees.

Here are the obtained results for tic-tac-toe endgame dataset:

```

1 Mean of the accuracy:  0.7888541666666666 --> 79%
2 Var of the accuracy:  6.982629681810922e-07
3 Best accuracy:  0.8958333333333334 --> 90%

```

Here are the obtained results for wine dataset:

```

1 Mean of the accuracy:  0.7822222222222223 --> 78%
2 Var of the accuracy:  1.2356500533455256e-05
3 Best accuracy:  1.0 --> 100%

```

Finally, confusion matrices were calculated for the tree which gives the best accuracy among 10-times-10-fold cross validation. It can be seen in Table II

TABLE II
CONFUSION MATRICES

Predicted Numbers	
Positive	Negative
59	5
5	27

(a) Confusion Matrix for tic-tac-toe endgame Dataset

Predicted Numbers		
Class "1"	Class "2"	Class "3"
6	0	0
0	8	0
0	0	4

(b) Confusion Matrix for wine Dataset

When the decision tree was built based on gain ratio instead of information gain, the performance of classifier for tic-tac-toe dataset was not changed significantly although there is a significant change for wine dataset. This behaviour is quite expected since generally gain ratio was used when a given input data has so many attribute values which force classifier based on information gain to choose that attribute. For tic-tac-toe dataset, all attributes has only three values: "x", "o" and "b", so building classifier based on gain ratio doesn't create any difference. This is also true for wine dataset, whose attributes are continious and they are converted to low and high than a certain threshold during the process, so they have only two different values. The observed difference was caused from the number of test samples for wine dataset. It is quite low, so small changes create big differences as percentage wise. In summary, obtained results are similar for both IG and GR case, since given attributes of the data don't have too many values.

III. UNDERSTANDING ATTRIBUTE AND CLASS NOISE ON THE PERFORMANCE OF ID3 CLASSIFIER

A. Analysis the Effect of Attribute Noise

Plots in Figure 1 show the classification accuracy respect to different attribute noises with the values of (5%, 10%, and 15%) for tic-tac-toe endgame and wine dataset. From these plots, we can conclude that corrupting test set is much more dangerous than corrupting training set since mean accuracy of 10-times-10-fold cross validation drops abruptly whenever test set is corrupted. Therefore, accuracy values are in the following order for both of the datasets:

Clean Training Set vs Clean Test Set (CvsC) > Dirty Training Set vs Clean Test Set (DvsC) > Clean Training Set vs Dirty Test Set (CvsD) > Dirty Training Set vs Dirty Test Set (DvsD)

The difference between datasets is reduction amount in accuracy is lower as noise level increases for wine dataset. It is not surprising since test set size is low for wine dataset, so corrupting a few samples make huge impact. One last note about the plots is that sometimes the trends in the curves are not as expected like CvsD case for wine dataset. The reason behind is the combination of randomness of the added noise and low test set size, i.e. sometimes flipping values randomly might not create too much change in the test set. For example, when test set has 8 samples, flipping randomly the values of 8 samples might create a new test set whose values are the same as old one like 6 out of 8 values might stay the same.

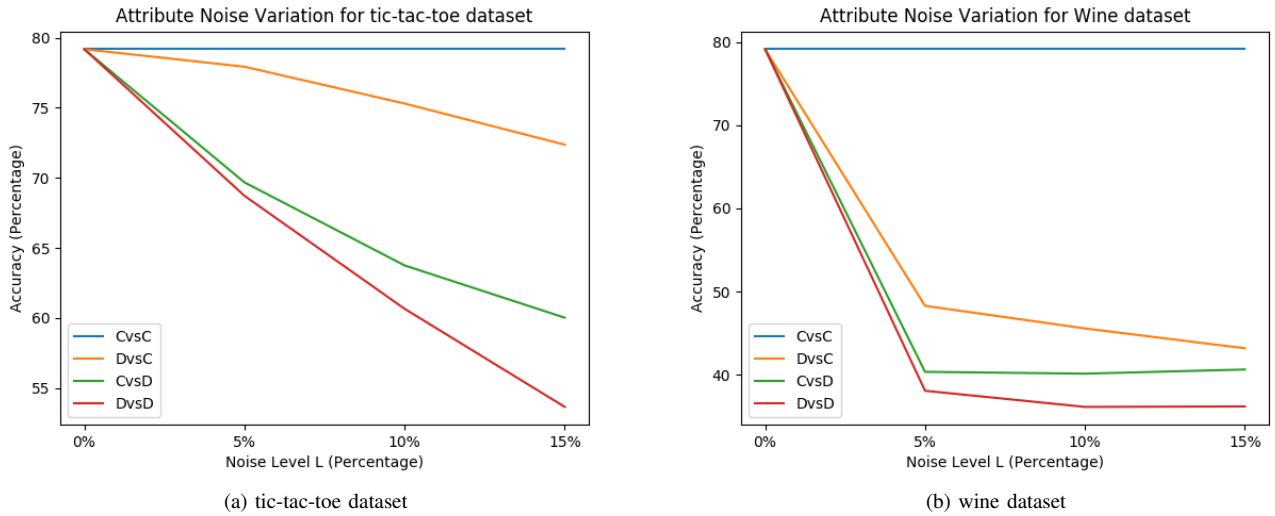


Fig. 1. Variation of Attribute Noise for Given Datasets

B. Analysis the Effect of Class-label Noise

Plots in Figure 2 show the classification accuracy respect to different class-label noises with the values of (5%, 10%, and 15%) for tic-tac-toe endgame and wine dataset. From these plots, we can infer that misclassifications in the training data affects the accuracy more than contradictory examples in the training data for both of the datasets. This makes sense since for contradictory examples, decision tree has both of the labels, so somehow it can self-correct itself or choose to not use these examples since same they point out different labels with same features. However, decision tree performs quite bad when there is misclassifications in the training data since it is not easy to detect misclassifications from clear training data.

The difference between these two datasets are the rate of reduction in accuracy. Reductions are much more dramatic in wine dataset due to small test sample size comparing to tic-tac-toe dataset.

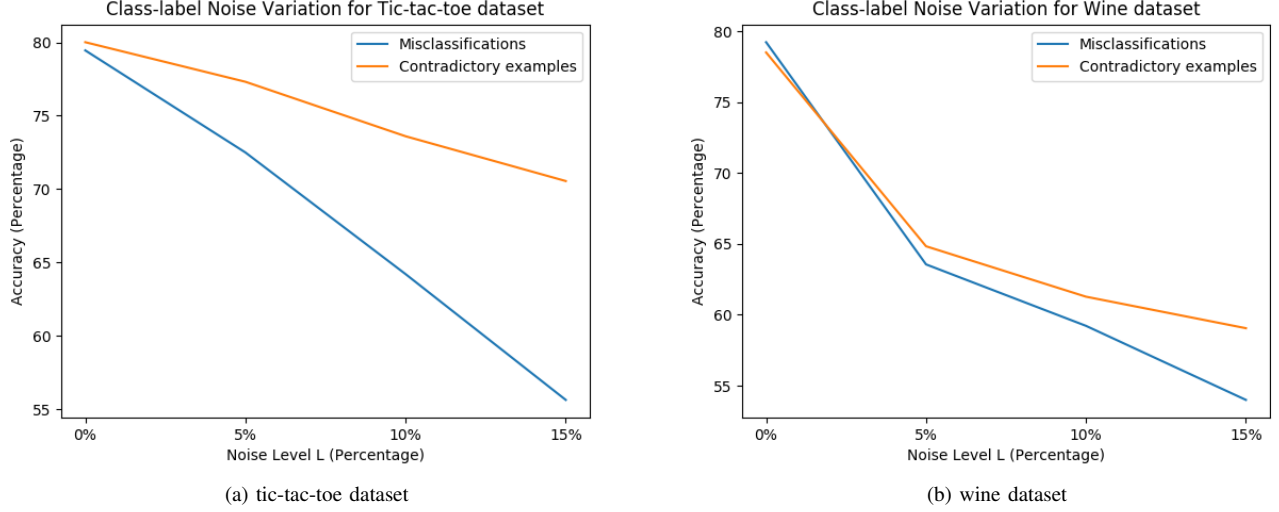


Fig. 2. Variation of Class-label Noise for Given Datasets

C. Prunning Variable and ID3

When there is not any pruning variable, having noise in the training data creates overfitting problem since ID3 tries to grow each branch of the decision tree in such a way that it can perfectly classify the data. When there is a noise in training data, the tree is erroneously fitting which creates overfitting problem or/and huge decision trees. This becomes quite crucial especially when the training data is too small like in our case with wine dataset. That's why post-pruning methods has been widely used in practice to overcome this problem.

D. Attribute Noise vs Class Noise

Class noise is more harmful than attribute noises since attribute noise changes the values of features while class noise changes directly the training labels. As a result, class noise also changes the balance of labels in the training data and it creates a class imbalance. For example, if you have a minority in your training data less say with label "1", if class noise changes these labels to majority, then you will lose quite precious information related to class "1" since there is limited number of examples. Therefore, in the end accuracy will drop significantly when this happens. It is problematic from decision tree perspective as well since it can change most common value in labels which will directly affect the tree structure negatively. Hence, class noise is usually more harmful than attribute noise. Attribute noise still can be managed to overcome if we have enough training samples since at the end decision tree itself is learning from given training data to map features with labels.