

# Final exam report of assignments

Subarna Khorshed

UG02-48-18-010

## 1 KNN:

The k-nearest neighbors (KNN) algorithm is a simple, supervised machine learning algorithm that can be used to solve both classification and regression problems. It's easy to implement and understand, but has a major drawback of becoming significantly slower as the size of that data in use grows.

### **Algorithm Explanation:**

- Take Data input from CSV file
- Decision tree is built as below-
  - o Find which attribute has the maximum information gain by finding the entropy for tuple.
  - o Now find the best split value of that particular attribute and save it in the TreeNode.
  - o Now as per the split separate the dataset into two parts - Left and Right and then recursively find the attribute with maximum gain with repeating the steps above.
  - o At the end of the recursion we will have the model built with leaf nodes representing the class variables - 'B', 'G', 'M', 'N'
- After the tree is built, we save the tree using pickle package in python
- To test the model on test dataset, we find confusion matrix and the accuracy of the model as below-
  - o Key concept is that we run our model using the `classify()` method which will traverse the model as per the each instance of the test dataset and finds out the predicted class.
  - o Accuracy is calculated using the number of instances

## 1.1 Code Image:

```
1  File: kNN.py
2  This file contains a single class definition from which
3  k-nearest neighbor models can be instantiated.
4  This class implements two methods: train and predict.
5  While the k-nearest neighbor algorithm does not require training,
6  the api has been defined like this to remain consistent with other
7  open-source machine learning libraries like scikit-learn.
8  """
9
10
11  import numpy as np
12  from src.distance import euclidean
13
14
15  class KNearestNeighbors:
16
17      def __init__(self, k=3, distance_metric=euclidean):
18          """Initialize k value and distance metric used for model."""
19          self.k = k
20          self.distance = distance_metric
21          self.data = None
22
23      def train(self, X, y):
24          """Zip labels and input data together for classification."""
25          # raise value error if inputs are wrong length or different types
26          if len(X) != len(y) or type(X) != type(y):
27              raise ValueError("X and y are incompatible.")
28          # convert ndarrays to lists
29          if type(X) == np.ndarray:
30              X, y = X.tolist(), y.tolist()
31          # set data attribute containing instances and labels
32          self.data = [X[i]+[y[i]] for i in range(len(X))]
33
34      def predict(self, a):
35          """Predict class based on k-nearest neighbors."""
36          neighbors = []
37          # create mapping from distance to instance
38          distances = {self.distance(x[:-1], a): x for x in self.data}
39          # collect classes of k instances with shortest distance
40          for key in sorted(distances.keys())[:self.k]:
41              neighbors.append(distances[key][-1])
42          # return most common vote
43          return max(set(neighbors), key=neighbors.count)
```

## 2 Decision Tree:

A decision tree is a decision support tool that uses a tree-like model of decisions and their possible consequences, including chance event outcomes, resource costs, and utility. It is one way to display an algorithm that only contains conditional control statements

### Algorithm Explanation

Decision Tree algorithm belongs to the family of supervised learning algorithms. Unlike other supervised learning algorithms, the decision tree algorithm can be used for solving regression and classification problems too. The goal of using a Decision Tree is to create a training model that can use to predict the class or value of the target variable by learning simple decision rules inferred from prior data. In Decision Trees, for predicting a class label for a record we start from the root of the tree. We compare the values of the root attribute with the records attribute. On the basis of comparison, we follow the branch corresponding to that value and jump to the next node.

## 2.1 Code Screenshot

```
9     def __init__(self,split,col_index):
10         self.col_id= col_index
11         self.split_value= split
12         self.parent=None
13         self.left= None
14         self.right= None
15
16     class Tree():
17
18         def __init__(self):
19             self.treemodel = None
20
21         def train(self,trainData):
22             ##Attributes/Last Column is class
23             self.createTree(trainData)
24
25
26         def createTree(self,trainData):
27             #create the tree
28             self.treemodel=build_tree(trainData,[])
29             saveTree(self.treemodel)
30
31         def accuracy_confusion_matrix(self,testData):
32             #prints the tree confusion matrix along with the accuracy
33             build_confusion_matrix(self.treemodel,testData)
34
35 #returns the best split on the data instance along
36 #with the splitted dataset and column index
37 def getBestSplit(data):
38     #set the max information gain
39     maxInfoGain = -float('inf')
40
41     #convert to array
42     dataArray = np.asarray(data)
43
44     #to extract rows and columns
45     dimension = np.shape(dataArray)
46
47     #iterate through the matrix
48     for col in range(dimension[1]-1):
49         dataArray = sorted(dataArray, key=lambda x: x[col])
50         for row in range(dimension[0]-1):
51             val1=dataArray[row][col]
52             val2=dataArray[row+1][col]
53             expectedSplit = (float(val1)+float(val2))/2.0
54             infoGain,l,r= calcInfoGain(data,col,expectedSplit)
55             if(infoGain>maxInfoGain):
56                 maxInfoGain=infoGain
57                 best= (col,expectedSplit,l,r)
58     return best
59
60 #This method is used to calculate the gain and returns
61 #the left and right data as per the split
62 def calcInfoGain(data,col,split):
63     totalLen = len(data)
64     infoGain = entropy(data)
65
66     left_data, right_data =splitDataColit(data,split,col)
```

### 3 Conclusion

I would like to thank to my teacher for giving me this amazing opportunity to learn these algorithms. I am sure this would be of help for my career of programming.