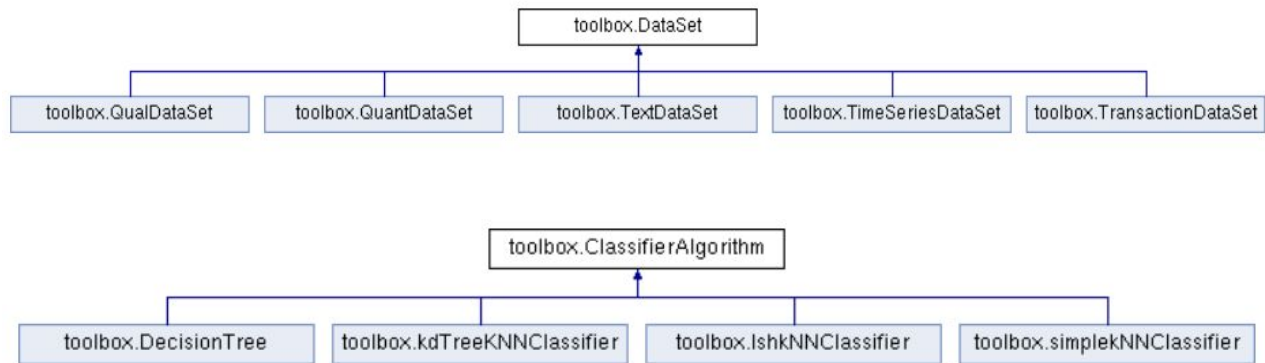


Summary of Object Oriented Toolbox Design

UML :



Summary :

Dataset: Super class contains common functions for all types of dataset, used to store, modify and return datasets.

TimeSeriesDataSet: Subclass inherits from Dataset that is used to store, modify and draw time series datasets. Includes time period selection and time type conversion functions.

TextDataSet: Subclass inherits from Dataset that is used to store, modify and explore text datasets. Includes stemming, lemmatization and stopwords removal to modify text.

QuantDataSet: Subclass inherits from Dataset that is used to store, modify and explore quantitative datasets. Explore the dataset with a line chart and a bar graph.

QualDataSet: Subclass inherits from Dataset that is used to store, modify and explore qualitative datasets. Explore the dataset with a pie chart and a bar graph.

TransactionDataSet: Subclass inherits from Dataset that is used to store, modify and explore transaction datasets. Contains inner class Rule for association rules calculation

ClassifierAlgorithm: Super class contains common functions for all classifiers, used to store training set, store train-test ratio and auto split the dataset to train and test by ratio.

simplekNNClassifier: Subclass inherits from ClassifierAlgorithm that uses k-nearest neighbors algorithm to classify the test set by training set, training label and stored k.

DecisionTree: Subclass inherits from ClassifierAlgorithm and Tree class to build a decision tree with training set and training label, pruning the tree and classifying the test set.

lshkNNClassifier: Subclass inherits from ClassifierAlgorithm using Locality Sensitive Hashing together with k-nearest neighbors algorithm to classify the test set.

kdTreeKNNClassifier: Subclass inherits from ClassifierAlgorithm to classify the test set by building kd Tree, an algorithm which uses a mixture of Decision trees and KNN.

Experiment: Experiment class that used to test compare between different algorithms. It contains functions for cross validation, confusion matrix and ROC curve.

HeterogenousDataSets: HeterogeneousDataSets class to store, load, clean, and explore a list of datasets. Contains functions for select and get type of datasets.

Summary of 3 advanced algorithms including time complexity analysis

Kd-KNN classifier Algorithm

Summary :

KD Tree is one such algorithm which uses a mixture of Decision trees and KNN to calculate the nearest neighbour. Theoretically, at each node we will save 3 things: dimension we split on, value we split on and tightest bounding box which contains all the points within that node. The algorithm could be divided into three parts. First, A Kd-tree class is built to save the information for classification. In the class, not only the distance is calculated, but also the tree is built. Then a KNN function is used to obtain the k-nearest neighbours. The second part is the node part, it saves information of the left child, the right child, the point and the axis. The third part is the large heap. Basically, We built a large list to save the heap and we make modifications on the heap. This class allows the action of adding new elements in, popping the top elements out and exchanging the elements to make adjustments to make sure that the heap is in order.

Actual code for createTree algorithm:

(This algorithm is under the class kdTreeKNNClassifier. A full code can be found in toolbox.)

```
def createTree(self, data, current_node, axis=0): # 1 step
    """
    function to create a kd-tree, recursion
    :type data: array-like, samples used to construct a kd-tree or sub_kd-t
    :type axis: int, between 0 and n_dim, dimension used to split data
    :type current_node: Node, the current 'root' node
    :return: None
    """
    if self.size == 0: # 2 steps: if, ==
        self.root = current_node # 1 step
        self.size += 1 # 1 step
    if data.shape[0] == 1: # if no more than one sample, then stop iterat:
        current_node.point = data[0, :] # 3 steps: data[0,:], current_no
        current_node.axis = axis # 2 steps: current_node.axis, =
        return # 1 step
    """
    step1: split the points with the median on this axis
    To find the median on target axis, we've got two ways:
    A. simply sort on target axis each time
    B. presort on each axis? but I haven't solved this at present
    """
    temp = data[data[:, axis].argsort()] # 4 steps: data[:,axis], data[:,
    med = int(len(temp)/2) if len(temp) % 2 == 0 else int((len(temp)-1)/2) # 2
    # find the 'first' med, this means that "<" goto left child, ">" goto r: # 9
    # 9 steps: if, len(), %, ==, <, >

    while med > 0 and temp[med,axis] == temp[med-1, axis]: # 7 steps: while
        med -= 1 # 1 step
    current_node.axis = axis # 2 steps: current_node.axis, =
    current_node.point = temp[med] # 2 steps: current_node.point, =
    tt = temp[med] # 2 steps: temp[med], tt=
    axis = (axis + 1) % self.n_dim # 4 steps: self.n_dim, %, axis+1
    if temp[:,med, :].shape[0] >= 1: # 4 steps: if, temp[:, temp[med].sh
        tt = temp[:,med, :] # 2 steps: temp[:, tt=
        current_node.left = self.Node() # 2+5 steps: 2 steps- current_no
        self.createTree(temp[:,med, :], current_node.left, axis) # T(n/2)+2
    if temp[:,(med+1), :].shape[0] >= 1: # 4 steps: if, temp[:, temp[med].sh
        tt = temp[:,(med+1), :] # 2 steps: temp[:, tt=
        current_node.right = self.Node() # 2+5 steps: 2 steps- current_no
        self.createTree(temp[:,(med+1), :], current_node.right, axis) # T(n/2)+2

    # Total:
    # T(n) = T(n/2)+70 = logn+70
    # T(n) is O(n) = logn
```

Total Time Complexity for kdTreeKNNClassifier:

(Details of time complexity calculation please refer to the code.)

p: len(self.testDS.shape[0]), n: length of testData

$$T(n) = (3+p)*\log n + 351p + 404$$

$$T(n) \text{ is } O(n) = p*\log n$$

Decision Trees Algorithm

Summary :

Decision Tree is the classification algorithm that is suitable for both continuous and discrete labels. By passing the data_type parameter as “continuous” or “discrete”, buildTree function will build the tree corresponding to the parameter with the corresponding functions. In the tree building procedure, each node created with 9 attributes, including the decided attributes position of the dataset, left and right node, decided value for discrete tree, threshold for continuous tree, probability for each label stored in the node, etc. The decided attribute and its threshold value or decided value for each node is calculated based on the minimal Gini index. After building the tree node by node, the pruning function can help to prune the tree by comparing the NH value stored in a node and the sum of NH value in its left and right nodes.

Actual code for buildTree algorithm:

(This algorithm is under the class DecisionTree. A full code can be found in toolbox.)

```
def buildTree(self, X, y, data_type='continuous', maxDepth=5, currentDepth = 0):
    # 1 step
    Gain_max, thre, d = 1, 0, 0
    attr = ()
    X_small, y_small, X_big, y_big = [], [], [], []
    # 3 steps: Gain_max, attr, d
    # 1 step
    # 4 steps: X_small, y_small, X_big, y_big
    # 3 steps: y.size,
    ## when there is more than one y label
    if y.size > 1:
        # 2 steps: if, ==
        if data_type == 'continuous':
            Gain_max, thre, d = self.attribute_based_on_Giniindex(X, y, data_type)
            # 46 dx + 22 ldx +
        elif data_type == 'discrete':
            Gain_max, attr, d = self.attribute_based_on_Giniindex(X, y, data_type)
            # 13 steps: if, or
            if Gain_max >= 0 and len(list(set(y))) > 1 and len(np.unique(X, axis=0)) > 1:
                # 2 steps: if, ==
                if data_type == 'continuous':
                    X_small, y_small, X_big, y_big = self.divide_group(X, y, thre)
                    # 18+4 steps: 18 s
                elif data_type == 'discrete':
                    X_small, y_small, X_big, y_big = self.divide_group_discrete(X, y, attr)
                    # 2 steps: +=1, ct
                    currentDepth+=1
                    left_branch = self.buildTree(X_small, y_small, data_type, maxDepth, currentDepth)
                    right_branch = self.buildTree(X_big, y_big, data_type, maxDepth, currentDepth)
                    # T(n/2)+1 steps,
                    # T(n/2)+1 steps,

    nh = self.NtHt(y)
    max_label = self.maxLabel(y)
    prediction_prob = self.predictScore(y)
    return self.decisionnode(d=d, attr=attr, thre=thre, NH=nh, lb=left_branch, rb=right_branch)

# 10 + 1 steps: 10 steps - se

# Total:
# x: len(inputted training set), d: len(X_attr), l: len(list(set(y_small))), k: len(label_count.keys()), t: len(list(total_label)), m: length of inputted training set, n: length of testData
# T(x,d,l,k,t) = 2T(n/2) + 46dx + 22 ldx - 22lm - 15x + 12 l + 36k + 18t + 148
# = logn + 46dx + 22 ldx - 22lm - 15x + 12 l + 36k + 18t + 148
# T(x,d,l,k,t) is O(x,d,l,k,t) = logn + ldx + k + t
```

Total Time Complexity for DecisionTree:

(Details of time complexity calculation please refer to the code.)

d: len(X_attr), l: len(list(set(y_small))), k: len(label_count.keys()),
t: len(list(total_label)), m: length of inputted training set, n: length of testData

$$T(d,l,k,t,m,n) = 3\log n + 46dm + 22 ldm - 22lm + 12 l + 36k + 18t + 6mn + 7m + 7n + 248$$

$$T(d,l,k,t,m,n) \text{ is } O(d,l,k,t,m,n) = \log n + ldm + k + t + mn$$

ROC Algorithm

Summary :

ROC is a two-dimensional graph where the true positive rate(TPR) lies on the y axis and false positive rate(FPR) is on the x axis. Before drawing the ROC curve, points representing the FPR and TPR are calculated by using different thresholds. Thresholds are usually scores of positive class for classification. For example, they could be prediction scores for positive classes. The scores are in-order, from the largest to the smallest. (Thresholds will be chosen from $+\infty$ to $-\infty$ according to the ranked scores.) Scores larger than the chosen threshold are classified as positive class, and scores less or equivalent to the chosen threshold are classified as negative class. Then TPR and FPR are calculated and stored. When there are instances with the same scores, there will be no extra point to be stored. Iterations used to calculate true positive rate and false positive rate starts from the first threshold equals the largest score until the last threshold equals the smallest score. For a two-class classification problem, ROC will choose a class as the positive class and the other as the negative class to evaluate the performance of classifiers. For a multiple-class classification problem, one class is chosen as a positive class, and the others are all classified as negative classes.

Actual code for ROC algorithm:

(This algorithm is under the class Experiment. A full code can be found in toolbox.)

[illegible]

Total Time Complexity for ROC:

(Details of time complexity calculation please refer to the code.)

m: length of inputted training set, n: length of testData, k: len(self.classifier)

$$T(m,n,k) = 9mk + 11k + 15n + 32$$

$T(m,n,k)$ is $O(m,n,k) = mk+n$