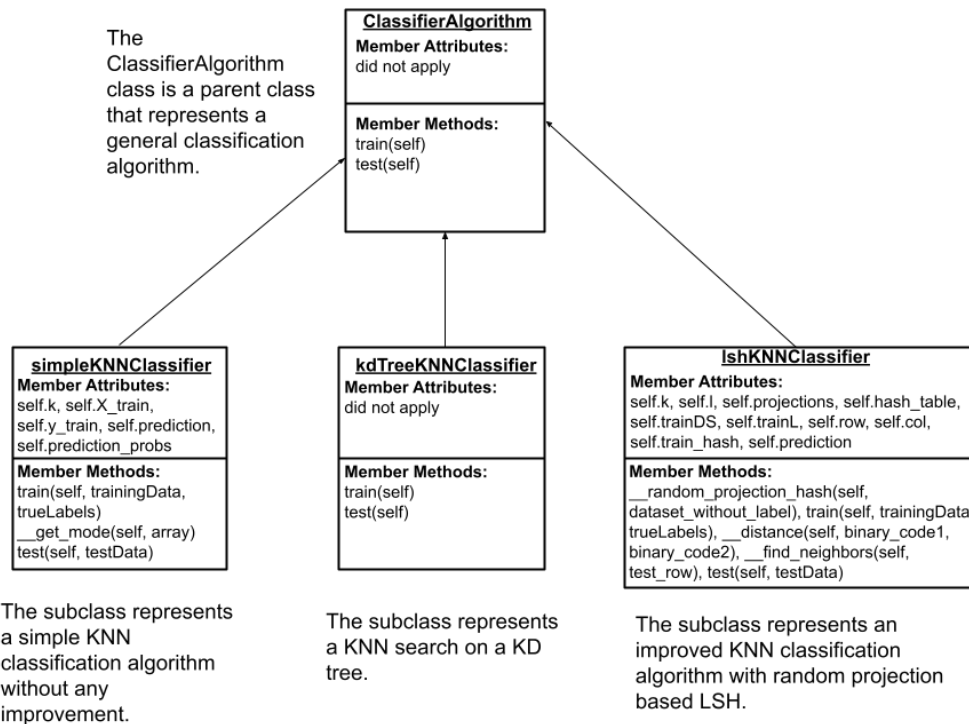
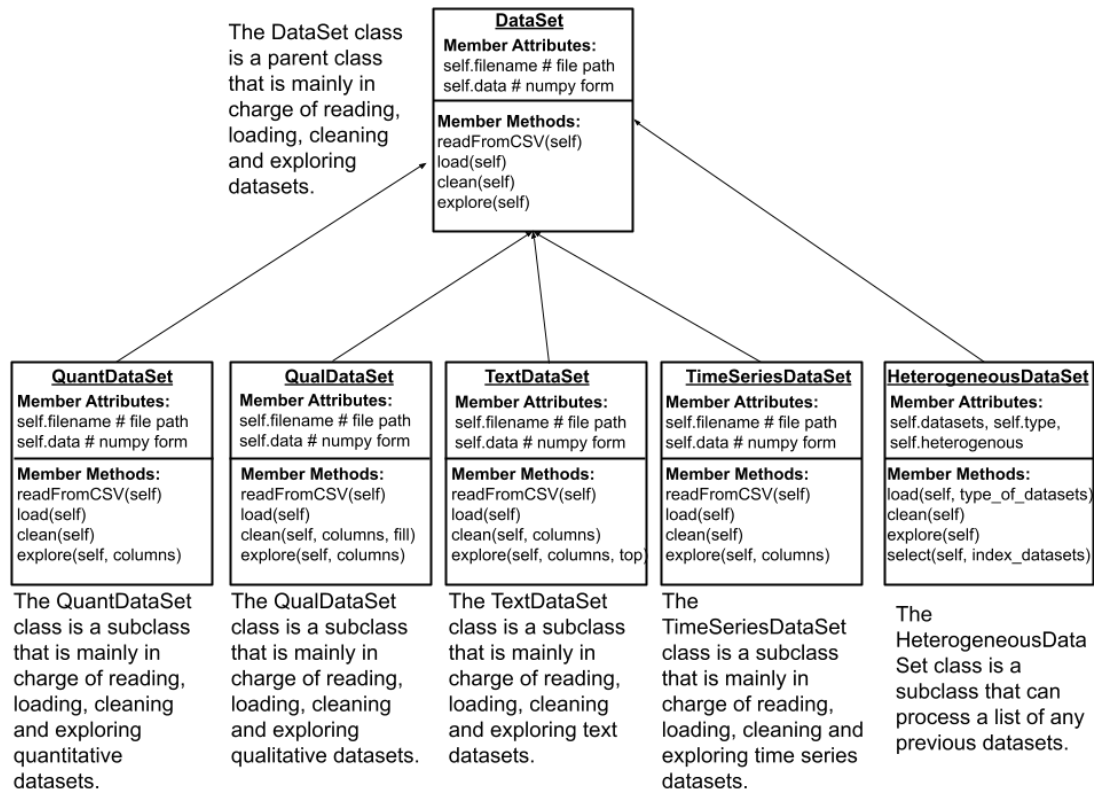


Summary of Object-oriented Toolbox Design



The Experiment class performs k-fold cross-validation and can evaluate the model performance with accuracy, confusion matrix, and ROC curves.

<u>Experiment</u>	
Member Attributes:	
self.data, self.labels, self.classifiers, self.predLabels, self.true, self.scores	
Member Methods:	
runCrossVal(self, k=5) score(self) __printScore(self) confusionMatrix(self) sortReverse(self, alist) ROC(self, prediction_probs, trueLabels)	

ROC Algorithm

1. Summary

The ROC algorithm is implemented according to the paper given. We can simply sort the probabilities reversely and update TP and FP for each predicted label. For each positive instance, we increase TP and for each negative instance, we increment FP. We also save each new ROC point to stack R. Finally, we plot the curve with the stack.

2. pseudo code or actual code

```
907 def ROC(self, prediction_probs, trueLabels):
908     """Produce a ROC plot.
909     Implemented according to the paper (An introduction to ROC analysis) by Tom Fawcett.
910
911     Keyword arguments:
912     prediction_probs -- a list of dictionaries containing prediction probabilities
913     trueLabels -- a list of true labels for the test set
914     """
915     legends = [] # step count # space count
916     labels_graph = np.unique(trueLabels) # 1 op # y
917                                           # n ops # y
918     for label_val in labels_graph: # y ops
919         probs = [d[label_val] if label_val in d.keys() else 0 for d in prediction_probs] # n ops # n
920
921         # sort the prediction probabilities in decreasing order
922         probs_sorted, idx_sorted = self.sortReverse(probs) # n*log(n) ops # 2n
923         # sort the test labels according to sorted indices
924         trueLabels = list(trueLabels) # 2 ops # n
925         labels_sorted = [trueLabels[i] for i in idx_sorted] # n ops # n
926
927         # initialize parameters
928         TP = 0 # 1 op # 1
929         FP = 0 # 1 op # 1
930         R = [] # a list of ROC points # 1 op # n
931         P = len([i for i in trueLabels if i == label_val]) # n ops # 1
932         N = len(trueLabels) - P # 3 ops # 1
933         f_prev = float('-inf') # 1 op # 1
934         i = 0 # 1 op # 1
935
936         # implement the algorithm
937         while i < len(labels_sorted): # n ops
938             f_i = probs_sorted[i] # 2 ops # 1
939             label = labels_sorted[i] # 2 ops # 1
940             if f_i != f_prev: # 1 op
941                 R.append((FP/N, TP/P)) # push point onto R # 1 op
942                 f_prev = f_i # 1 op
943             if label == label_val: # 1 op
944                 TP += 1 # true positive # 1 op
945             else: # 1 op
946                 FP += 1 # false positive # 1 op
947             i += 1 # 1 op
948
949         val_1 = (FP/N, TP/P) # 1 op # 1
950         R.append(val_1) # push (1, 1) onto R # 1 op
951         legends.append(f'Class {label_val}') # 1 op
952         plt.plot(*zip(*R)) # 1 op
953
954         plt.plot([0, 1], [0, 1], 'k--') # 1 op
955         plt.xlabel('False Positive Rate') # 1 op
956         plt.ylabel('True Positive Rate') # 1 op
957         plt.title('ROC') # 1 op
958         plt.legend(legends) # 1 op
959         plt.show() # 1 op
960
961     # n: length of test data
962     # y: number of unique labels
963     # y is much less than n. Thus, y can be regarded as a constant when computing time complexity.
964
965     # T(n) = 1 + n + y * (nlogn + 3n + 12n + 14) + 6
966     # Tight-fit upperbound for T(n): O(n * log(n))
967
968     # S(n) = 6n + 2y + 9
969     # Tight-fit upperbound for S(n): O(n)
```

3. time complexity analysis: shown above.

kd-trees Algorithm

1. Summary

To build a KNN-KD tree, first, select its root node and select an axis. Then sort the first dimensional data and found the median. Then divide the data into root nodes, the left subtrees, and the right subtrees. The above process goes recursively until all data is divided. The search process includes a forward search and a retrospective search. If the data to be searched is smaller than the axis data, it will go to the left subtree. Otherwise, it will take the right subtree until the last point is the leaf node. Then update the points repeatedly in the search path until we find the closest point.

2. pseudo code or actual code

Algorithm 2: KNN-KD-tree building (KNNKDTB)

Input: P , a set of training points; $depth$, the current depth

Output: $treeroot$, the root of the KD-tree storing P

1. $Treeroot \leftarrow root(P)$
 2. While ($P \neq null$)
 3. select axis = $depth \bmod k$
 4. sort P and select median by axis from P
 5. $root.location \leftarrow median$
 6. $root.leftChild \leftarrow KNNKDTB(\text{points in } P \text{ before median, } depth+1)$
 7. $root.rightChild \leftarrow KNNKDTB(\text{points in } P \text{ after median, } depth+1)$
 8. End while
 9. Output $treeroot$
-

Algorithm 3: K Nearest-neighbor Classification based on KD-tree (KNN-KD-tree)

Input: $treeroot$, the root of a KD-tree; $test_point$, a point of test data

Output: $predict_label$, the label of the point

1. While (root is not a leaf)
 2. $searchPath.add(root)$
 3. $root[axis] > test_point[axis] ? Searchtree(rightChild) :$
 $Searchtree(leftChild)$
 4. End while
 5. For point in $searchPath$:
 6. $nearest_dis \leftarrow Compute_distance (point, test_point)$
 7. If ($|point[axis] - test_point[axis]| > |test_point[axis] - root[axis]|$)
 8. Travel ($root.nextchild$)
 9. $dis \leftarrow Compute_distance (childpoint, test_point)$
 10. If ($dis < nearest_dis$)
 11. $nearest_dis \leftarrow dis$
 12. End if
 13. End for
 14. End for
 15. $predict_label \leftarrow nearestPoint.label$
 16. Output $predict_label$
-

3. time complexity analysis

The time complexity of a KNN search on a KD-tree is $O(n^{1-(1/k)} + m)$, where m is the number of nodes and k is the dimension of the KD-tree. In high-dimensional data, it is not as efficient as the simple KNN search.

Lsh-knn Algorithm

1. Summary

First, we create an $m \times l$ random matrix A from the standard normal distribution. Then we project the training data onto l dimensional space with this random matrix A. Then convert the $n \times l$ dimensional matrix to n l -bits binary codes. For each point in the test set, we can compute the Hamming distance between it and each training point. Finally, we sort the distances and get the k neighbors.

2. pseudo code or actual code (helper functions not include)

```
637 def train(self, trainingData, trueLabels):
638     """Return the hash indexes of the training data.
639
640     Keyword arguments:
641     trainingData -- data to train the classifier
642     trueLabels -- training labels
643     """
644     self.trainDS = trainingData           # step count      # space count
645     self.trainL = trueLabels              # 1 op             # n
646     self.row, self.col = self.trainDS.shape # 1 op             # n
647     self.train_hash = self.__random_projection_hash(self.trainDS) # 2 ops            # n + m
648     return self.train_hash                # n*m*l+3n+ml+2 ops # n
649
650 # Dimension of training data: n * m
651 # l: length of binary code
652
653 # T(n) = n*m*l + 3n + ml + 7
654 # Tight-fit upperbound: O(n * m * l)
655
656 # S(n) = 4n + m
657 # Tight-fit upperbound: O(n + m)
```

```
692 def test(self, testData):
693     """Return the prediction result and scores.
694
695     Keyword arguments:
696     testData -- data to test the classifier
697     """
698     # store the prediction result
699     self.prediction = []           # step count      # space count
700     # transform the test data
701     testData_hash = self.__random_projection_hash(testData) # 1 op             # n
702     # for each test row
703     for row in testData_hash:      # n ops
704         # get the nearest k neighbors
705         neighbor = self.__find_neighbors(row) # nlogn+4nl+4n+k+4 ops # k
706         # count the labels
707         lab, count = np.unique(neighbor, return_counts=True) # klogk ops        # 2k
708         # get the mode
709         predicted_class = lab[np.argmax(count)] # k ops            # k
710         # append the predicted label to the result list
711         self.prediction.append(predicted_class) # 1 op
712     # return predicted labels
713     return self.prediction         # 1 op
714
715 # Dimension of test data: n * m
716 # l: length of binary code
717 # k: number of nearest neighbors
718
719 # T(n) = n^2 * log(n) + 4n^2 * l + 4n^2 + nml + ml + n*k*log(k) + 2nk + 8n + 4
720 # Tight-fit upperbound: O(n^2 * log(n))
721
722 # S(n) = 2n + 4k
723 # Tight-fit upperbound: O(n + k)
724
725 # For my implementation, the time complexity does not improve compared to the simple KNN algorithm.
726 # Mainly due to the inefficient sorting method. In theory, the LSH KNN algorithm should have an
727 # improved time complexity compared to the simple KNN algorithm.
728
729 # The accuracy of LSH KNN is less than the accuracy of simple KNN. This is acceptable since LSH
730 # algorithms often require lots of tuning. Perhaps the current parameters still need to be adjusted
731 # to get the optimal result.
```

3. time complexity analysis: shown above.