

```

1 '''it returns an array of objects containing group and
2 feature list'''
3 def getDataFromFile(fileName):
4     data = []
5     file_in = open(fileName, 'r')
6
7     # read every line
8     for instance in file_in.readlines():
9         features = [] # create a feature list
10        obj = {} # create an object to store group and
11        feature information
12        instance = instance.strip() # remove leading and
13        trailing spaces
14
15        values = instance.split(' ') # split row based on
16        spaces
17        obj['group'] = float(values[0]) # store first value
18        as group
19
20        # read the remaining line and stores every feature
21        in features list
22        for feature in values[1:]:
23            if(feature == ''):# omit white space
24                continue
25            features.append(float(feature))
26
27        obj['features'] = features # add feature list to
28        the object
29        data.append(obj)
30
31    return data
32
33
34    '''Normalization using Z-Score'''
35    import statistics
36    def normalization(data):
37        # making deep copy of data
38        normalizedData = []
39        for instance in data:
40            tempObject = {}
41            tempObject['group'] = instance.get('group')
42            tempObject['features'] = [feature for feature in
43            instance.get('features')]
44            normalizedData.append(tempObject)
45
46
47        # loop through each features and normalize it
48        for i in range(len(normalizedData[0].get('features'))):
49            featureList = []
50
51            for instance in normalizedData:
52                featureList.append(instance.get('features')[i])

```

```

43
44     mean = statistics.mean(featureList)
45     stdev = statistics.stdev(featureList)
46
47     for instance in normalizedData:
48         feature = instance.get('features')[i]
49         feature = (feature - mean) / stdev
50         instance['features'][i] = feature
51
52     return normalizedData
53
54 import math
55 '''it accepts array of objects as training data and just
list of features for test data'''
56 def nearestNeighbourClassifier(trainingDataSet,
testInstance):
57     distances = []
58     for trainingInstance in trainingDataSet: #get object
from array
59         distanceBetweenFeature = float(0)
60         # loop through each feature and calculate distance
61         for trainingFeature, testingFeature in zip(
trainingInstance.get('features'), testInstance):
62             distanceBetweenFeature += (trainingFeature -
testingFeature)**2
63
64         # store value of euclidean distance with group
number as tuples
65         euclideanDistance = math.sqrt(
distanceBetweenFeature)
66         distances.append((euclideanDistance,
trainingInstance.get('group')))
67
68     # sort distances and returns the group number
69     distances = sorted(distances)
70     return distances[0][1]
71
72
73 def leaveOneOutValidator(data):
74     correctGuess = 0
75     # select each index to leave as test data
76     for leaveIndex in range(len(data)):
77         testData = data[leaveIndex]
78         if(leaveIndex == 0):
79             trainingDataSet = data[1:]
80         else:
81             # add the remaining instances to the training
set
82             firstSubset = data[0:leaveIndex]
83             secondSubset = data[leaveIndex+1:]

```

```

84         trainingDataSet = firstSubset + secondSubset
85
86         # call nearest neighbour algo for each set and
87         # check if it has detected the group correctly
88         guessedGroup = nearestNeighbourClassifier(
89             trainingDataSet, testData.get('features'))
90         if(guessedGroup == testData.get('group')):
91             correctGuess = correctGuess+1
92
93     # returns the overall percentage
94     validationScore = correctGuess/len(data)
95     return validationScore * 100
96
97     '''it returns a set of features for a given indices'''
98 def extractFeatures(data, featureIndices):
99     newDataSet = []
100    for i in range(len(data)):
101        obj = {}
102        extractedFeatures = []
103        for index in featureIndices:
104            extractedFeatures.append(data[i].get('features')
105                )[index])
106        obj['group'] = data[i].get('group')
107        obj['features'] = extractedFeatures
108        newDataSet.append(obj)
109
110    return newDataSet
111
112 def forwardSelection(data):
113     print('Beginning Search.\n')
114     featureCount = len(data[0].get('features'))
115
116     # indices of features which improves the model
117     indices = []
118     bestScores = [] # best scores on each iteration
119
120     checkedLocalMaxima = False
121     # maximum number of iteration
122     for i in range(featureCount):
123         best = 0.0
124         newIndex = 0
125
126         # add new feature on each iteration for evaluation
127         for j in range(featureCount):
128             if(j in indices):
129                 continue
130
131             # copy previously selected indices
132             temp = [val for val in indices]

```

```

131          # add new index
132          if(len(temp) != 0):
133              if(j < temp[0]):
134                  temp.insert(0, j)
135              else:
136                  temp.append(j)
137          else:
138              temp.append(j)
139
140          # extract subset of features and groups
141          # according to new list of features to evaluate
142          newDataSet = extractFeatures(data, temp)
143          accuracyRate = leaveOneOutValidator(newDataSet)
144
145          # update accuracy rate on every addition of
146          # features
147          print('\t Using feature(s)' + str([index+1 for
148          index in temp]) + ' accuracy is ' + str(accuracyRate))
149          if(accuracyRate > best):
150              best = accuracyRate
151              newIndex = j
152
153
154          # for local maxima
155          if(len(bestScores) != 0 and best < bestScores[len(
156          bestScores) - 1][0]):
157              if(checkedLocalMaxima):
158                  print('\nAddition of features is not
159                  improving the model\n')
160              else:
161                  print('\n(Warning, Accuracy has decreased
162                  ! Continuing search in case of local maxima)')
163              checkedLocalMaxima = True
164
165
166          tempBestScoreIndices = [index+1 for index in
167          indices]
168          if(len(tempBestScoreIndices) != featureCount):
169              print('\nFeature(s) set' + str(
170              tempBestScoreIndices) + ' was best, accuracy is ' + str(
best) + '\n')
171          else:
172              print()

```

```

171     bestScores.append((best, tempBestScoreIndices))
172     bestScores = sorted(bestScores)
173
174
175     print('Finished Search!! The best feature subset is '
176     + str(bestScores[len(bestScores) - 1][1]) + ', which has
177     an accuracy of '+str(bestScores[len(bestScores) - 1][0]))
178
179
180 def backwardElimination(data):
181     print('Beginning Search.\n')
182     featureCount = len(data[0].get('features'))
183
184     indices = [i for i in range(featureCount)] # indices
185     of features which improves the model
186
187     score0fAllFeatures = leaveOneOutValidator(data) # accuracy when all features are selected
188
189     bestScores = score0fAllFeatures
190     checkedLocalMaxima = False
191
192     # maximum number of iteration
193     for i in range(featureCount):
194         scoreList = []
195         accuracyImproved = False
196
197         # remove a feature on each iteration
198         for j in range(len(indices)):
199             temp = [index for index in indices] # copy
200             previously selected index list
201             temp.pop(j) # remove a feature
202
203             # get features based on the indices of new
204             features
205             newDataSet = extractFeatures(data, temp)
206             accuracyRate = leaveOneOutValidator(newDataSet)
207
208             # evaluate accuracy and update if needed
209             print('\t Using feature(s)' + str([index + 1
210             for index in temp]) + ' accuracy is ' + str(accuracyRate))
211             if (accuracyRate >= bestScores):
212                 bestScores = accuracyRate
213                 accuracyImproved = True
214
215             scoreList.append((accuracyRate, temp))
216
217             # check for local maxima
218             if(not accuracyImproved):

```

```

213             if(checkedLocalMaxima):
214                 print('\nAddition of features is not
improving the model\n')
215                 break
216             else:
217                 print('\n(Warning, Accuracy has decreased
! Continuing search in case of local maxima)')
218                 checkedLocalMaxima = True
219
220             scoreList = sorted(scoreList)
221             indices = scoreList[len(scoreList) - 1][1]
222             print('\nFeature(s) set' + str([index + 1 for
index in indices]) + ' was best, accuracy is ' + str(
bestScores) + '\n')
223
224             print('Finished Search!! The best feature subset is '
+ str([index + 1 for index in indices]) + ', which has an
accuracy of ' + str(bestScores))
225
226
227
228 def main():
229     print('Welcome to Calvin Ng \'s Feature Selection
Algorithm')
230     fileName = input('Type in the name of the file to test
: ')
231
232     print('\nType the number of the algorithm you want to
run.\n')
233     print('1)Forward Selection')
234     print('2)Backward Elimination')
235
236     algoType = input('\t\t\t')
237
238     data = getDataFromFile(fileName)
239
240     print('This dataset has '+str(len(data[0]).get(
'features'))+' features (not including the class attribute
), with '+str(len(data))+ ' instances.')
241
242     print('\nPlease wait while I normalize the data...', end
=' ')
243     normalizedData = normalization(data)
244     print('Done!')
245
246     accuracyOfAllFeatures = leaveOneOutValidator(
normalizedData)
247     print('\nRunning nearest neighbor with all '+str(len(
data[0].get('features')))+ ' features, using "leaving-one-
out" evaluation, I get an accuracy of '+str(

```

```
247 accuracyOfAllFeatures)+"%\n")  
248  
249     if(algoType == '1'):  
250         forwardSelection(normalizedData)  
251     elif(algoType == '2'):  
252         backwardElimination(normalizedData)  
253     else:  
254         print('Please input correct type')  
255  
256  
257 if __name__ == '__main__':  
258     main()
```