

САНКТ-ПЕТЕРБУРГСКИЙ ПОЛИТЕХНИЧЕСКИЙ
УНИВЕРСИТЕТ им. ПЕТРА ВЕЛИКОГО

Институт прикладной математики и механики

Высшая школа прикладной математики и вычислительной физики

Отчёт
по лабораторной работе №7 «Алгоритмы композиции»
по дисциплине «Системы искусственного интеллекта»

Студентка гр. 3630201/70101 _____ О. В. Саксина

Преподаватель _____ Л. В. Уткин

Содержание

1	Задание 1	3
1.1	Постановка задачи	3
1.2	Реализация	3
2	Задание 2	3
2.1	Постановка задачи	3
2.2	Реализация	3
3	Задание 3	4
3.1	Постановка задачи	4
3.2	Реализация	4
	Приложение	5

1 Задание 1

1.1 Постановка задачи

Исследуйте зависимость тестовой ошибки от количества деревьев в ансамбле для алгоритма adaboost.M1 на наборе данных Vehicle из пакета mlbench (обучающая выборка должна состоять из 7/10 всех прецедентов, содержащихся в данном наборе данных). Постройте график зависимости тестовой ошибки при числе деревьев, равном 1, 11, 21, . . . , 301, объясните полученные результаты.

1.2 Реализация

Результат представлен на рис. 1. При увеличении количества деревьев в ансамбле до 80, тестовая ошибка уменьшается. Затем колеблется в пределах от 0.225 до 0.3, принимая наименьшее значение на 100 деревьях.

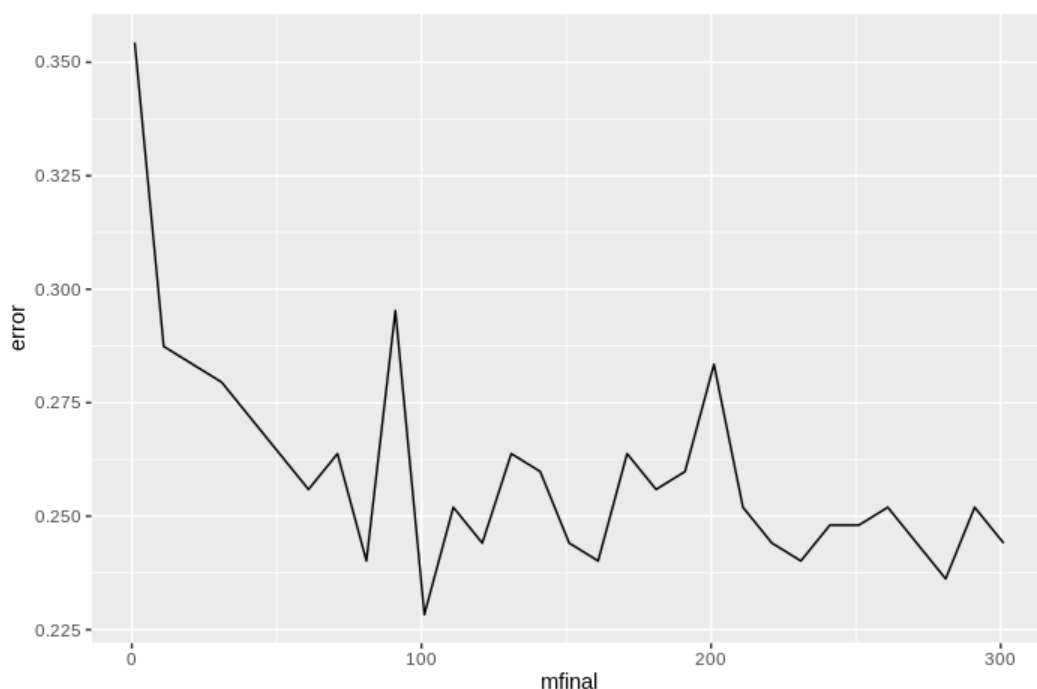


Рис. 1: Зависимость ошибки от количества деревьев в ансамбле для алгоритма adaboost

2 Задание 2

2.1 Постановка задачи

Исследуйте зависимость тестовой ошибки от количества деревьев в ансамбле для алгоритма bagging на наборе данных Glass из пакета mlbench (обучающая выборка должна состоять из 7/10 всех прецедентов, содержащихся в данном наборе данных). Постройте график зависимости тестовой ошибки при числе деревьев, равном 1, 11, 21, . . . , 201, объясните полученные результаты.

2.2 Реализация

Результат представлен на рис. 2. При увеличении количества деревьев в ансамбле до 30 тестовая ошибка уменьшается. Затем колеблется в пределах от 0.24 до 0.29, принимая наименьшее значение на 70 деревьях.

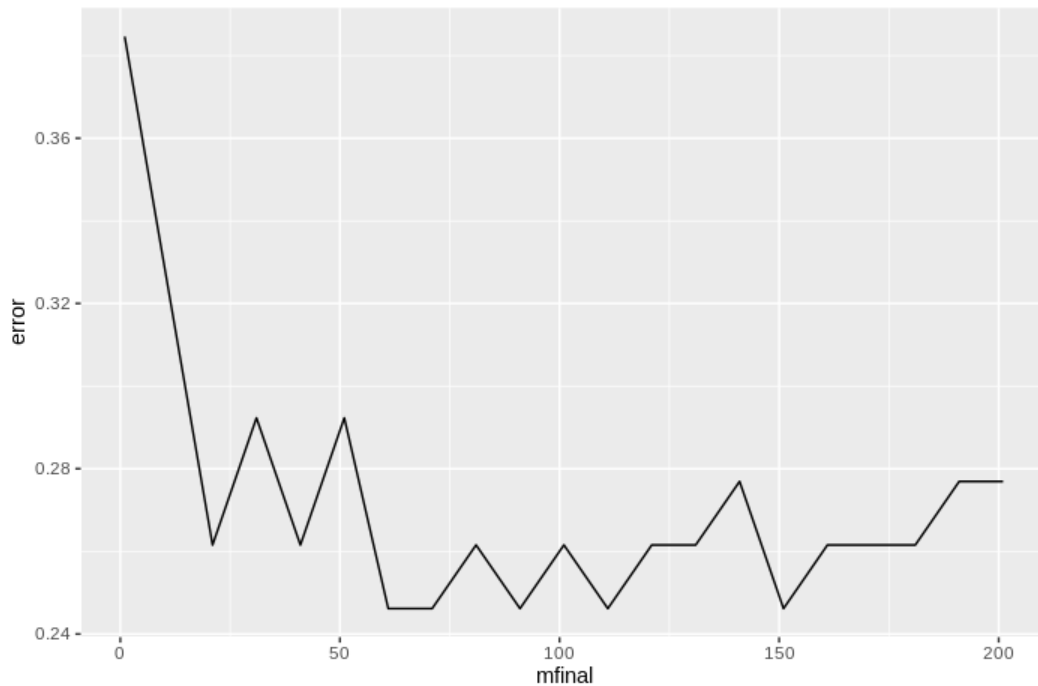


Рис. 2: Зависимость ошибки от количества деревьев в ансамбле для алгоритма bagging

3 Задание 3

3.1 Постановка задачи

Реализуйте бустинг алгоритм с классификатором K ближайших соседей. Сравните тестовую ошибку, полученную с использованием данного классификатора на наборах данных Vehicle и Glass, с тестовой ошибкой, полученной с использованием единичного дерева классификации.

3.2 Реализация

Описание алгоритма на естественном языке:

1. Присвоить каждому элементу из набора данных вес, равный 1.
2. Разделить данные на тестовые и обучающие.
3. Обучающие данные методом leave one out настроить следующим образом:
 - (a) Найти k ближайших соседей для i -го элемента.
 - (b) Определить количество голосов за каждый класс, сложив по отдельности веса соседей каждого класса. Присвоить i -му элементу класс, за который отдано больше всех голосов.
 - (c) Если этот класс не совпадает с реальным классом i -го элемента, то для каждого соседа из k ближайших соседей увеличить на 0.5 вес тех, кто имеет реальный класс i -го элемента, и уменьшить на 0.5 вес тех, кто имеет назначенный класс i -го элемента.
 - (d) Повторить шаги а-с 2 раза
4. Предсказать классы элементов из тестового набора путём поиска k ближайших соседей в обучающих данных с изменёнными весами.

Средняя тестовая ошибка данного классификатора на наборе данных Glass при $k = 8$ — 0.35. Средняя тестовая ошибка единичного дерева классификации глубины 4 — 0.31.

Средняя тестовая ошибка данного классификатора на наборе данных Vehicle при $k = 8$ — 0.41. Средняя тестовая ошибка единичного дерева классификации глубины 4 — 0.32.

Приложение

Задание 1

```
1 library(ggplot2)
2 library(rpart)
3 library(mlbench)
4 library(adabag)
5
6 data(Vehicle)
7 l = length(Vehicle[,1])
8 sub = sample(1:l, 7*l/10)
9 mfina1 = seq(1, 301, 10)
10 maxdepth = 5
11 res = rep(0, length(mfina1))
12 for(i in 1:length(mfina1)){
13   Vehicle.adaboost = boosting(Class ~., data=Vehicle[sub,], mfina1=mfina1[i], maxdepth=maxdepth)
14   Vehicle.adaboost.pred = predict.boosting(Vehicle.adaboost, newdata=Vehicle[-sub, ])
15   res[i] = Vehicle.adaboost.pred$error
16 }
17 df = data.frame(mfina1 = mfina1, error = res)
18 ggplot(df, aes(x = mfina1, y = error)) +
19   geom_line()
```

Задание 2

```
1 library(ggplot2)
2 library(rpart)
3 library(mlbench)
4 library(adabag)
5
6 data(Glass)
7 l = length(Glass[,1])
8 sub = sample(1:l, 7*l/10)
9 mfina1 = seq(1, 201, 10)
10 maxdepth = 5
11 res = rep(0, length(mfina1))
12 for(i in 1:length(mfina1)){
13   Glass.adaboost = bagging(Type ~., data=Glass[sub,], mfina1=mfina1[i], maxdepth=maxdepth)
14   Glass.adaboost.pred = predict.bagging(Glass.adaboost, newdata=Glass[-sub, ])
15   res[i] = Glass.adaboost.pred$error
16 }
17 df = data.frame(mfina1 = mfina1, error = res)
18 ggplot(df, aes(x = mfina1, y = error)) +
19   geom_line()
```

Задание 3

```
1 import numpy as np
2 from collections import Counter
3 import csv
4 import random
5 import copy
6 from sklearn.tree import DecisionTreeClassifier
7
8 def distance(instance1, instance2):
9     instance1 = np.array(instance1)
10    instance2 = np.array(instance2)
11    return np.linalg.norm(instance1 - instance2)
12
13 def get_distances(tr, te):
14     train = copy.deepcopy(tr)
15     test = copy.deepcopy(te)
16     for i in range(len(te)):
```

```

17     test[i].pop(-1)
18     for i in range(len(tr)):
19         train[i].pop(-1)
20     distances = {}
21     test_ind = 0
22     for i in test:
23         distances[test_ind] = []
24         train_ind = 0
25         for j in train:
26             distances[test_ind].append([train_ind, distance(i, j)])
27             train_ind += 1
28         test_ind += 1
29     return distances
30
31 def get_neighbors(train, train_labels, test_num, k, distances):
32     #neighbour = (index, dist, class, weight)
33     neighbors = []
34     distances[test_num].sort(key=lambda x: x[1])
35     neighbors = copy.deepcopy(distances[test_num][:k])
36     for i in range(k):
37         neighbors[i].append(train_labels[neighbors[i][0]])
38         neighbors[i].append(train[neighbors[i][0]][-1])
39     return(neighbors)
40
41 def vote(neighbors):
42     class_counter = Counter()
43     for neighbor in neighbors:
44         class_counter[neighbor[2]] += neighbor[3]
45     return class_counter.most_common(1)[0][0]
46
47 def add_weight(test_num, right_label, wrong_label, train, train_labels, neighbors):
48     for n in neighbors:
49         ind = n[0]
50         if train_labels[ind] == right_label:
51             train[ind][-1] += 0.5
52         elif train_labels[ind] == wrong_label:
53             train[ind][-1] -= 0.5
54
55 def read_data(file):
56     data = []
57     with open(file) as csvfile:
58         reader = csv.reader(csvfile)
59         for row in reader:
60             data.append(row)
61     for row in range(1, len(data)):
62         for el in range(len(data[row]) - 1):
63             data[row][el] = float(data[row][el])
64         try:
65             data[row][-1] = int(float(data[row][-1]))
66         except:
67             pass
68     return data
69
70 def knn(train0, main_train_labels):
71     for s in range(len(train0)):
72         print('\n-----')
73         print(str(s) + '-----\n')
74         k = 6
75         train_labels = copy.deepcopy(main_train_labels)
76         test = [train0[s]]
77         test_labels = [train_labels[s]]
78         if s == 0:
79             train = train0[1:]
80         elif s == len(train0):
81             train = train0[:s]
82         else:

```

```

83     train = train0[:s] + train0[s+1:]
84     d = get_distances(train, test)
85     votes = []
86     for test_ind in range(len(test)):
87         neighbors = get_neighbors(train, train_labels, test_ind, k, d)
88         votes.append(vote(neighbors))
89     wrongs = {}
90     for i in range(len(votes)):
91         if votes[i] != test_labels[i]:
92             wrongs[i] = (test_labels[i], votes[i])
93             print(str(i) + ' ') predicted: ' + str(votes[i]) + ' real: ' + str(test_labels[i]))
94     for i in wrongs:
95         add_weight(i, wrongs[i][0], wrongs[i][1], train, train_labels, get_neighbors(train, train_labels, i, k+1, d))
96
97     print('\n')
98     if wrongs:
99         k += 1
100        votes = []
101        for test_ind in range(len(test)):
102            neighbors = get_neighbors(train, train_labels, test_ind, k, d)
103            votes.append(vote(neighbors))
104        wrongs = {}
105        for i in range(len(votes)):
106            if votes[i] != test_labels[i]:
107                wrongs[i] = (test_labels[i], votes[i])
108                print(str(i) + ' ') predicted: ' + str(votes[i]) + ' real: ' + str(test_labels[i]) )
109        for i in wrongs:
110            add_weight(i, wrongs[i][0], wrongs[i][1], train, train_labels, get_neighbors(train, train_labels, i, k+1, d))
111        print('\n')
112
113        if wrongs:
114            k += 1
115            votes = []
116            for test_ind in range(len(test)):
117                neighbors = get_neighbors(train, train_labels, test_ind, k, d)
118                votes.append(vote(neighbors))
119            wrongs = {}
120            for i in range(len(votes)):
121                if votes[i] != test_labels[i]:
122                    wrongs[i] = (test_labels[i], votes[i])
123                    print(str(i) + ' ') predicted: ' + str(votes[i]) + ' real: ' + str(test_labels[i]) )
124            for i in wrongs:
125                add_weight(i, wrongs[i][0], wrongs[i][1], train, train_labels, get_neighbors(train, train_labels, i, k+1, d))
126    return train0
127
128 def main():
129     #data = read_data('glass.csv')
130     data = read_data('vehicle_csv.csv')
131     rows = len(data)
132     data[0].insert(-1, "weight")
133     for i in range(rows):
134         data[i].insert(-1, 1)
135     cols = len(data[0])
136     names = data[0]
137     data.pop(0)
138     random.shuffle(data)
139     train_size = round(0.7 * rows)
140     train0 = data[:train_size]
141     test0 = data[train_size:]
142     main_train_labels = [row[-1] for row in train0]
143     main_test_labels = [row[-1] for row in test0]
144     for row in test0:
145         row.pop(-1)
146     for row in train0:
147         row.pop(-1)
148     #tree

```

```

149     clf = DecisionTreeClassifier(max_depth = 4, random_state = 0)
150     clf.fit(train0, main_train_labels)
151     clf.predict(test0)
152     score = clf.score(test0, main_test_labels)
153     #knn
154     train = knn(train0, main_train_labels)
155     k = 8
156     test = test0
157     train_labels = copy.deepcopy(main_train_labels)
158     test_labels = copy.deepcopy(main_test_labels)
159     d = get_distances(train, test)
160     votes = []
161     for test_ind in range(len(test)):
162         neighbors = get_neighbors(train, train_labels, test_ind, k, d)
163         votes.append(vote(neighbors))
164     wrongs = {}
165     print('\n test \n')
166     for i in range(len(votes)):
167         if votes[i] != test_labels[i]:
168             wrongs[i] = test_labels[i]
169             print(str(i) + ') predicted: ' + str(votes[i]) + ' real: ' + str(test_labels[i]))
170     print('\n knn missclasification = ' + str(len(wrongs)/len(test)))
171     print('\n tree missclasification = ' + str(1 - score))
172
173
174 main()

```
