

# Bios/CS 534 Project 2

Chenxi Cai

March 28, 2018

## 1 Problem 1

With perceptron function, the error rate is

0.033333

(1)

Python codes of problem 1:

```
1 import numpy as np
2 import pandas as pd
3
4 #define functions of prediction
5 def predict(x,w,y):
6     '''
7     :param x: one row of data samples
8     :param w: data weight
9     :param y: one row of data labels
10    :return: predict label
11    '''
12    y_p = np.dot(x,w)*y
13    return y_p
14 #define function of perceptron
15 def perceptron(X,w,Y,lr):
16     '''
17     :param X: data samples
18     :param w: data weight
19     :param Y: data labels
20     :return: weight vector as a numpy array
21     '''
22     flag = True
23     while(flag):
24         total_error = 0
25         for i, x in enumerate(X):
26             x = np.insert(x,2,1) # Add bias = 1 in x
27             y_p = predict(x,w,Y[i])
28             if (y_p <= 0):
29                 total_error += 1
30                 w = w + lr*x*Y[i] #new weight,lr = learning rate
31         if (total_error == 0):
32             flag = False #loop until all training data are predict correct
33     return w
34
35 #prediction of testing data and error rate
36 def Error_rate(X2,w,Y2):
37     '''
38     :param X2: Testing data samples
39     :param w: data weight
40     :param Y: Testing data labels
41     :return: error_rate
42     '''
43     s = Y2.shape[0]
44     n = 0
```

```

45     for i, x in enumerate(X2):
46         x2 = np.insert(x,2,1)
47         y2_p = predict(x2,w,Y2[i])
48         if (y2_p <= 0):#error_rate = #points whose predict class is different
            form original one/#total points
49             n += 1
50     return n/float(s)
51
52 #read files
53 f = pd.read_csv('/Users/ccai28/Desktop/hw2_data_1.txt',sep='\t', header = None
    , skiprows = 1)
54
55 #read training data
56 X1 = f.loc[:69,[0,1]].values
57 Y1 = f.loc[:69,[2]].values
58
59 #read testing data
60 X2 = f.loc[70:,[0,1]].values
61 Y2 = f.loc[70:,[2]].values
62
63 #set learning rate and initial weight
64 lrate = 1
65 w = np.ones(3)
66
67 #new weight after training
68 w = perceptron(X,w,Y,lrate)
69
70 #Error rate
71 error_rate = Error_rate(X2,w,Y2)
72
73 print('perceptron: \nerror rate = {:.f} \n'.format(error_rate))

```

## 2 Problem 2

With a function of Adaboost, the error rate is shown in Table below

Iteration	3	5	10	20
Error rate	0.167	0.167	0.133	0.133

Plot of Error rate vs iteration is shown as Fig. 1

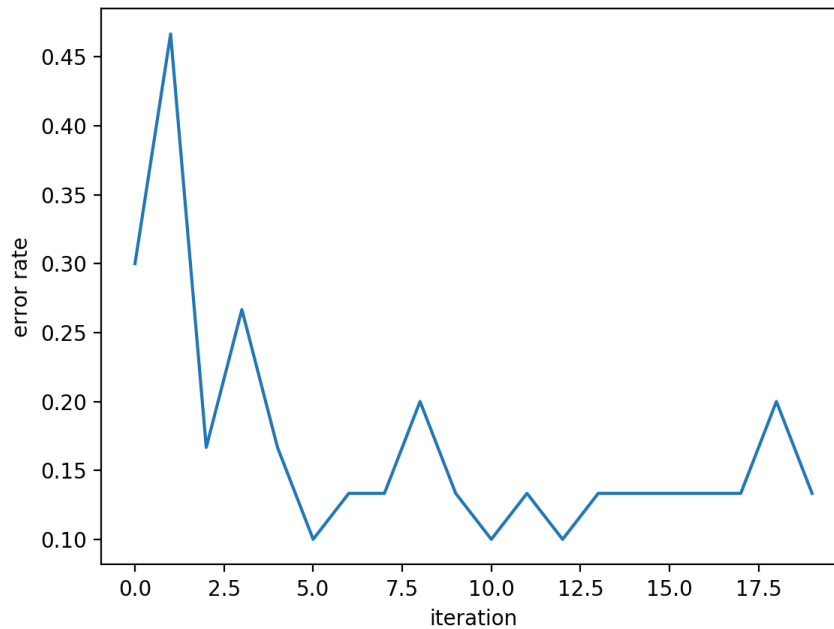


Figure 1: Plots of Adaboost error rate

### Python codes of problem 2:

```

1 import numpy as np
2 import pandas as pd
3 import math
4 import matplotlib.pyplot as plt
5
6 #define functions of classifier
7 def lineclassify(value,x,label):
8     '''
9     :param value: double,the standard value to compare with
10    :param x: training data,double,data axis value
11    :label: 0,left; 1,right
12    :left : x <= value: predict y is -1; x > value: predict y is 1
13    :right: x <= value: predict y is 1; x > value: predict y is -1
14    :return: predict y,int(-1 or 1)
15    '''
16    if (label == 0):
17        if (x <= value):
18            y-p = -1
19        else:
20            y-p = 1
21    elif (label == 1):
22        if (x <= value):
23            y-p = 1
24        else:
25            y-p = -1
26
27    return y-p
28
29
30 #Classify all training data with one weak classifier (xij)(loop all training
    data with specific xij)
31 def Weighterror(value,X,Y,W,label):
32     '''

```

```

33 :param value: double,the standard value to compare with
34 :param X: training data,list (n*1),data set in one axis (eg.X[:,0])
35 :param Y: training data,list (n*1),all data labels
36 :param W: list (n*1),weight set of data
37 :label: 0,left; 1,right
38 :left : x <= value: predict y is -1; x > value: predict y is 1
39 :right: x <= value: predict y is 1; x > value: predict y is -1
40 :return: double, weighterror = sum(misclassified data's value on one axis
    * weight)
    '''
41
42 weighterror = 0
43 for i, x in enumerate(X):
44     y_p = lineclassify(value,x,label)
45     y_dot = y_p*Y[i]
46
47     if (y_dot <= 0):
48         weighterror += W[i]
49
50 return weighterror
51
52 #The sum of weight
53 def Sumweight(W):
54     '''
55     :param W: list (n*1),weight set of data (all equal to 1/N),N is number of
    training data
56     :return: double,sum of weight
    '''
57
58     sumweight = 0
59     for i, w in enumerate(W):
60         sumweight += W[i]
61     return sumweight
62
63 #new weight after iteration
64 def changeweight(value,X_line,Y,W,label,a):
65     '''
66     :param value: double,the best value to compare with
67     :param X: training data,list (n*1),data set in one axis (eg.X[:,0])
68     :param Y: training data,list (n*1),all data labels
69     :param W: list (n*1),weight set of data
70     :para a: alpha wi = wi*exp(a)
71     :label: 0,left; 1,right
72     :left : x <= value: predict y is -1; x > value: predict y is 1
73     :right: x <= value: predict y is 1; x > value: predict y is -1
74     :return: list,new weight
    '''
75
76     for i,x in enumerate(X_line):
77         y_p = lineclassify(value,X_line[i],label)
78         y_dot = y_p*Y[i]
79
80         if (y_dot <= 0):
81             # print W[j],alpha[i]
82             W[i] = W[i]*np.exp(a)#calculate new wi
83
84         else:
85             continue
86     return W
87
88 #Find the best xij to classify traning data (loop for every possible xij)
89 def Buildline(X,Y,W):
90     '''
91     :param X: training data,list (n*2),data set in all axis (X[xi,x2])
92     :param Y: training data,list (n*1),all data labels
93     :param W: list (n*1),weight set of data
94     :return: list,best xij with smallest weighterror and its label in l/r, X1/
    X2
95
96         x_best = [best value,[l(0)/r(1),X1(0)/X2(1)]
97     :label: left:0 ; right:1; X1:0 ; X2:1
    '''

```

```

98 weight_x1_l = []
99 weight_x1_r = []
100 weight_x2_l = []
101 weight_x2_r = []
102 weight_x = []
103 x_best = []
104 X_list = []
105
106 X1 = X[:,0] #all x1
107 X2 = X[:,1] #all x2
108
109 X_list.append(X1)
110 X_list.append(X1)
111 X_list.append(X2)
112 X_list.append(X2) #X_list = [[X1],[X1],[X2],[X2]]
113
114
115 for i, x1 in enumerate(X1):
116     weighterror1_l = Weighterror(X1[i],X1,Y,W,0) #weighterror of X1,left
117     weight_x1_l.append(weighterror1_l)
118
119     weighterror1_r = Weighterror(X1[i],X1,Y,W,1)
120     weight_x1_r.append(weighterror1_r) #weighterror of X1,right
121
122 weight_x.append(weight_x1_l) #put (l,x1) in list , index = 0
123 weight_x.append(weight_x1_r) #put (r,x1) in list , index = 1
124
125
126 for i, x2 in enumerate(X2):
127     weighterror2_l = Weighterror(X2[i],X2,Y,W,0) #weighterror of X2,left
128     weight_x2_l.append(weighterror2_l)
129     weighterror2_r = Weighterror(X2[i],X2,Y,W,1) #weighterror of X2,right
130     weight_x2_r.append(weighterror2_r)
131
132
133 weight_x.append(weight_x2_l) #put (l,x2) in list , index = 2
134 weight_x.append(weight_x2_r) #put (r,x2) in list , index = 3
135 #weight_x = [[weight_x1_l], [weight_x1_r], [
weight_x2_l], [weight_x2_r]]
136
137 weight_x_min_list = np.min(weight_x, axis=1) #minimum of every (l/r,x1/x2)
list
138 #weight_x_min_list = [min(
weight_x1_l),min(weight_x1_r),min(weight_x2_l),min(weight_x2_r)]
139
140 weight_x_min = min(weight_x_min_list)#minimum in all list , smallest
weighterror
141
142 index_list = np.argmin(weight_x_min_list) #index of four lists who has the
smallest weighterror
143 #index_list = the index in
weight_x_min_list
144 # = the index in weight_x = best
value with smallest weighterror in which list , (l/r,x1/x2),between 0 and 3
145
146 label_list = [[0,0],[1,0],[0,1],[1,1]] #build label list to represent all
(l/r,x1/x2)
147 #label_list = [[l,X1],[r,X1],[l,X2
],[r,X2]]
148 label = label_list[index_list] #find the label(0 or 1) of smallest
weighterror
149
150 index = weight_x[index_list].index(weight_x_min)#index of best value,0-69
151
152 best_value = X_list[index_list][index]
153
154 x_best.append(best_value)
155 x_best.append(label) #x_best = [best value , [l(0)/r(1),X1(0)/X2(1)]

```

```

156                                     #best value = x_best[0], 1/r = x_best[1][0], X1/X2 =
    x_best[1][1]
157
158     return x_best
159
160 #Build Adaboost training
161 def Adaboost(iteration,X,Y,W,X_test,Y_test):
162     '''
163     :param iteration: times loop weak lineclassify
164     :param X: list (n*2),training data set in all axis (X[xi,x2])
165     :param Y: list (n*1),training data, all data labels
166     :param X_test: list (n*2),testing data set in all axis (X[xi,x2])
167     :param Y_test: list (n*1),testing data, all data labels
168     :param W: list (n*1),initial weight set of data (all equal to 1/N),N is
    number of training data
169     :return: list, error rate, contains every iteration error rate
170     '''
171     error_rate = []
172
173     Y_testp = np.zeros(Y_test.shape[0])
174     for i in range(iteration):
175         n = 0
176         x_best = Buildline(X,Y,W)
177
178         value, label_index, line_index = x_best[0], x_best[1][0], x_best[1][1]
179         #value, label_index, line_index = best value, 1/r, X1/X2
180         X_line = X[:,line_index] #x1 best: X_line = X[:,0]; x2 best: X_line = X
181        [:,1];
182
183         weighterror = Weighterror(value,X_line,Y,W,label_index)
184
185         sumweight = Sumweight(W)
186         error = weighterror/float(sumweight) #calculate error
187         print("errorm:")
188         print(error)
189
190         #print weighterror,sumweight,error
191         a = np.log((1-error)/float(error)) #calculate alpha
192         #print error,a,value_best[i]
193
194         W = changeweight(value,X_line,Y,W,label_index,a) #new weight
195
196         s_test = X_test.shape[0] #number of data in X_test
197         for j in range(s_test):
198             x_test = X_test[:,line_index][j]
199             Y_testp[j] += AdaClassify(a,x_best,x_test)
200
201             y_dot = Y_testp[j] * Y_test[j]
202             if (y_dot <= 0):
203                 n += 1
204         error_rate.append(n/float(s_test)) #error rate of every iteration
205         print(error_rate[i])
206
207     return error_rate
208
209 #Adaboost classify(combination of weak classifier)
210 def AdaClassify(alpha,value_best,x_test):
211     '''
212     :param alpha: to build adaboostc classifier: sum(alpha[i]*Gi(x))
213     :param value_best: list, = [[best value1, [1(0)/r(1), X1(0)/X2(1)], [best
    value2, [1(0)/r(1), X1(0)/X2(1)]]...]
214     :param x_test: element, testing data of one axis
215     :return: predict y with alpha
216     '''
217     value, label_index, line_index = value_best[0], value_best[1][0], value_best
218     [1][1] #value, label_index, line_index = best value, 1/r, X1/X2
219     y_testp = lineclassify(value,x_test,label_index)*alpha #adaboost
220     classifier

```

```

217
218     return y_testp
219
220 #read files
221 f = pd.read_csv('/Users/ccai28/Desktop/hw2_data_1.txt', sep='\t', header = None
222               , skiprows = 1)
223
224 #read training data
225 X = f.loc[:69,[0,1]].values
226 Y = f.loc[:69,[2]].values
227
228 #read testing data
229 X_test = f.loc[70:,[0,1]].values
230 Y_test = f.loc[70:,[2]].values
231
232 #number of training data
233 s = X.shape[0]
234
235 #set initial weight
236 W = np.full((s,1),1/float(s))
237
238 #iteration
239 iteration = 20
240
241 error_rate = Adaboost(iteration,X,Y,W,X_test,Y_test)
242 print error_rate[iteration-1]
243
244 x_plot = list(range(iteration))
245
246 plt.plot(x_plot, error_rate)
247 plt.xlabel('iteration')
248 plt.ylabel('error rate')
249 plt.show()

```

### 3 Problem 3

By using SVM in sklearn, the best parameter and error rate of three kernel methods are shown in Table below

	Radical kernel	Sigmoid kernel	Polynomial kernel
Best parameter (gamma/degree)	0.02	0.06	2
Error rate	0.443	0.493	0.426

Plots of Scores vs gamma or degree of three kernels are shown as Fig. 2, 3, 4

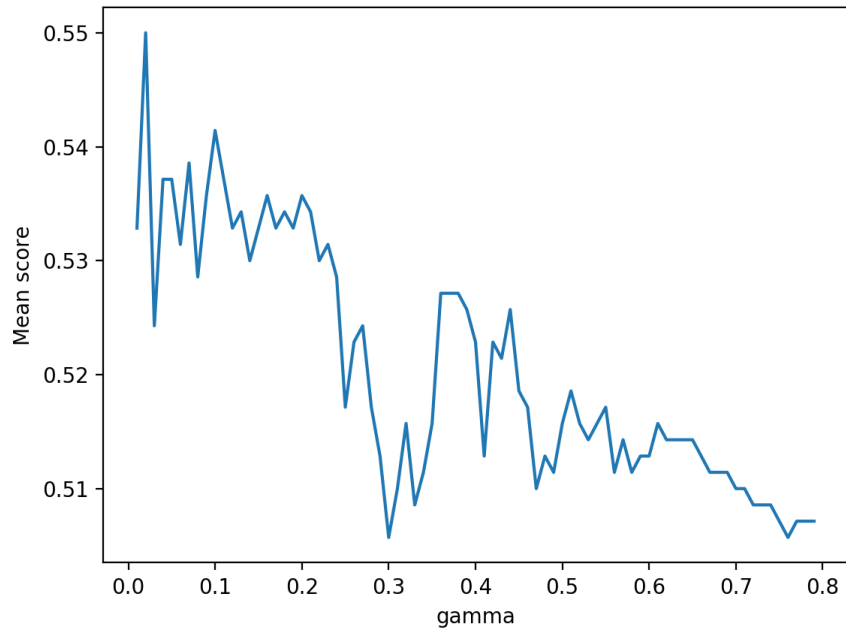


Figure 2: Plots of radical kernel

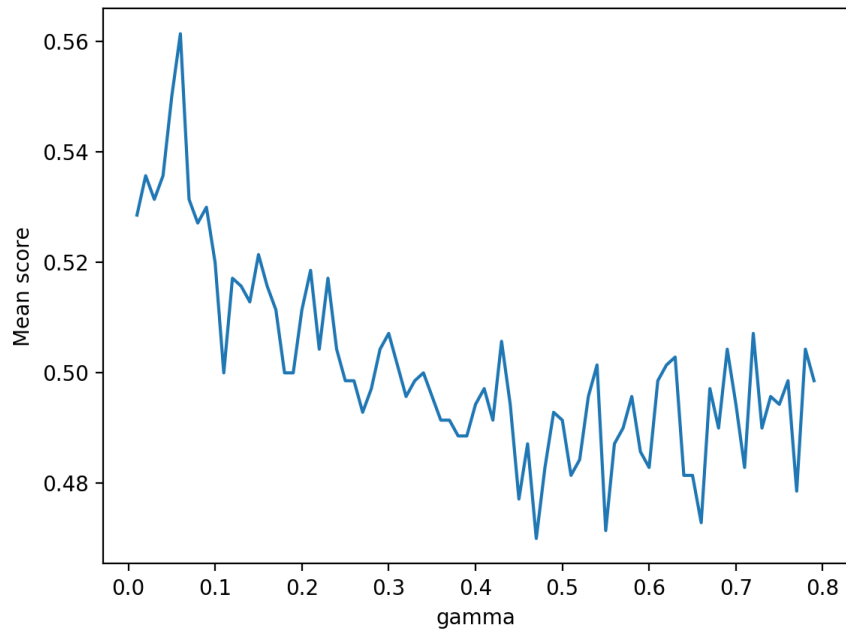


Figure 3: Plots of sigmoid kernel



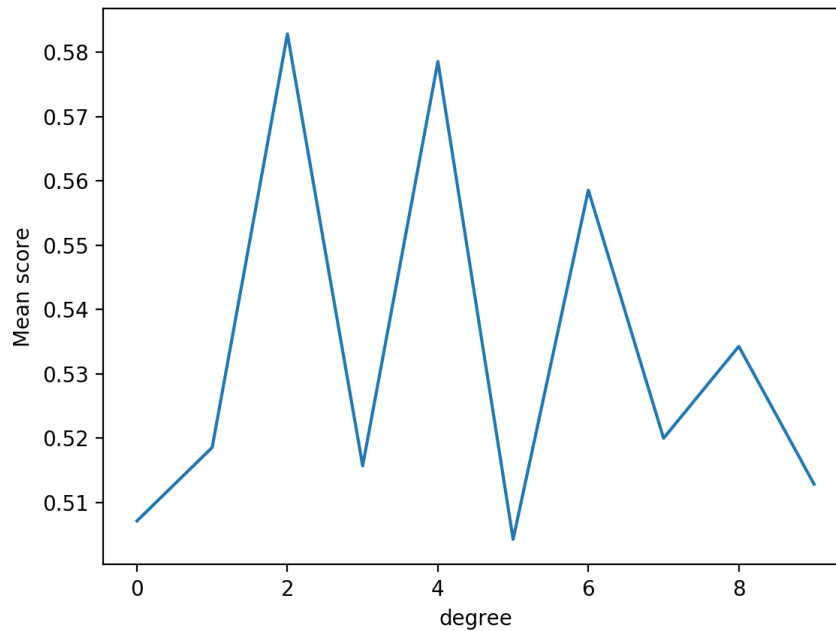


Figure 4: Plots of polynomial kernel

### Python codes of problem 3:

```

1 import numpy as np
2 import pandas as pd
3 import matplotlib.pyplot as plt
4
5 from sklearn.svm import SVC
6 from sklearn.grid_search import GridSearchCV
7
8 def SVM(param_grid, cv, kernel):
9     """
10     :param param_grid: parameter in SVM grid
11     :param cv : cv-fold
12     :param kernel: string, the name of kernel
13     :return: svm grid
14     """
15     grid = GridSearchCV(SVC(kernel=kernel), param_grid, cv = cv)
16     return grid
17
18 def Best_estimator(grid, X, Y):
19     """
20     :param grid: SVM grid
21     :param X : Training data X
22     :param Y : Training data Y
23     :return: best estimator after fitting training data. It has smallest error
24             and best parameter (gamma or degree)
25     """
26     grid.fit(X, Y)
27     best_estimator = grid.best_estimator_
28     print("The best classifier is: ", best_estimator)
29     return best_estimator
30
31 def plot(grid, C_range, para_range, s):
32     """
33     :param grid: SVM grid

```

```

33 :param C_range : range of C parameter
34 :param para_range : range of gamma or degree parameter
35 :param s : string, name of parameter('gamma' or 'degree')
36 :return: plot of score vs parameter
37 '''
38 scores = [x[1] for x in grid.grid_scores_]
39
40 #plot
41 plt.plot(para_range, scores)
42
43 plt.legend()
44 plt.xlabel(s)
45 plt.ylabel('Mean score')
46 plt.show()
47
48 def Error_rate(X,Y,X2,Y2,best_estimator):
49     '''
50     :param best_estimator: best_estimator with best parameter
51     :param X: list (n*20), training data set in all axis (X[xi,x2])
52     :param Y: list (n*1), training data, all data labels
53     :param X2: list (n*20), testing data set in all axis (X[xi,x2])
54     :param Y2: list (n*1), testing data, all data labels
55     :return: error-rate
56     '''
57     s = X2.shape[0]
58     n = 0
59     best_estimator.fit(X,Y)
60     Y_p = best_estimator.predict(X2)
61     for i, y in enumerate(Y_p):
62         if (y != Y2[i]):#error_rate = #points whose predict class is different
63             n += 1
64     return n/float(s)
65
66
67 #read files
68 f = pd.read_csv('/Users/ccai28/Desktop/hw2_data_2.txt',sep='\t', header = None
69               , skiprows = 1)
70
71 #read training data
72 X = f.loc[:699,:19].values
73 Y = f.loc[:699,20].values
74
75 #read testing data
76 X_test = f.loc[700,:19].values
77 Y_test = f.loc[700,20].values
78
79 #set parameter for Gridsearch
80 C_range = [1]
81 gamma_range = np.arange(0.01, 0.8,0.01)
82 degree_range = np.arange(0, 10,1)
83
84 param_grid_g = dict(gamma=gamma_range)
85 param_grid_d = dict(degree=degree_range)
86 cv = 10
87
88 #radical kernel
89 grid_r = SVM(param_grid_g,cv,'rbf')
90 best_estimator_r = Best_estimator(grid_r,X,Y)
91 plot(grid_r,C_range,gamma_range,'gamma')
92
93 error_rate_r = Error_rate(X,Y,X_test,Y_test,best_estimator_r)
94 print('radical kernel: \nerror rate = {:.f} \n'.format(error_rate_r))
95
96
97 #sigmoid kernel
98 grid_s = SVM(param_grid_g,cv,'sigmoid')

```

```

99 best_estimator_s = Best_estimator(grid_s,X,Y)
100 plot(grid_s,C_range,gamma_range,'gamma')
101
102 error_rate_s = Error_rate(X,Y,X_test,Y_test,best_estimator_s)
103 print('sigmoid kernel: \nerror rate = {f} \n'.format(error_rate_s))
104
105 # polynomial kernel
106 grid_p = SVM(param_grid_d,cv,'poly')
107 best_estimator_p = Best_estimator(grid_p,X,Y)
108 plot(grid_p,C_range,degree_range,'degree')
109
110 error_rate_p = Error_rate(X,Y,X_test,Y_test,best_estimator_p)
111 print('polynomial kernel: \nerror rate = {f} \n'.format(error_rate_p))

```

## 4 Problem 4

The plot of OOB error rate and number of trees are shown in Fig. 5,

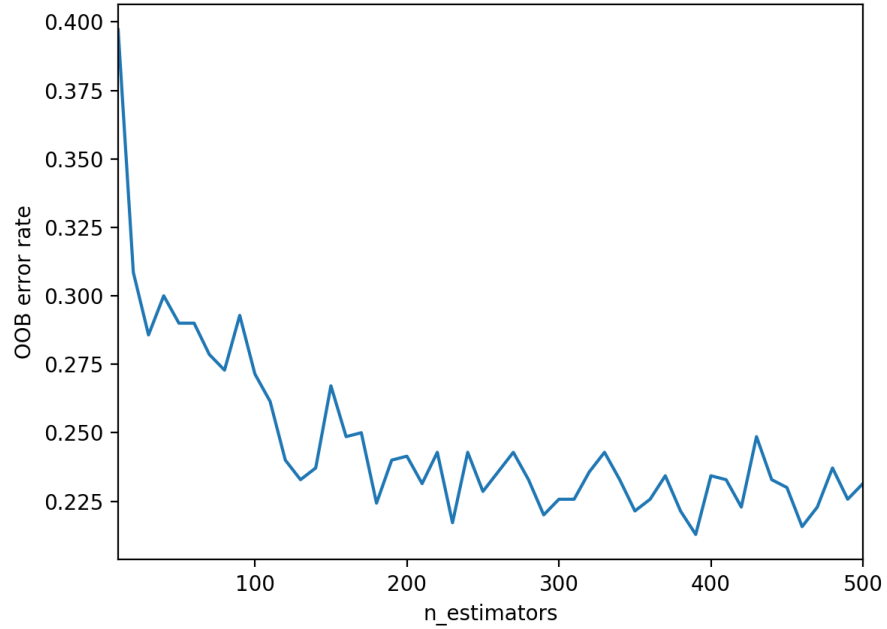


Figure 5: plots of OOB error rate vs number of trees

The selected number of trees is 210, and the error rate is

$$0.223333 \quad (2)$$

The feature ranking is:

rank	index of variables	importances
1	4	(0.129395)
2	9	(0.119601)
3	14	(0.101843)
4	0	(0.041882)
5	10	(0.041192)
6	2	(0.041054)
7	12	(0.039885)
8	1	(0.039104)
9	16	(0.038682)
10	3	(0.038664)
11	11	(0.038254)
12	7	(0.038123)
13	5	(0.037836)
14	19	(0.037576)
15	18	(0.037312)
16	8	(0.036664)
17	6	(0.036404)
18	13	(0.036342)
19	15	(0.035886)
20	17	(0.034302)

#### Python codes of problem 4:

```

1 import pandas as pd
2 import numpy as np
3 import matplotlib.pyplot as plt
4
5 from sklearn.ensemble import RandomForestClassifier
6
7
8 def Error_rate(X2,Y2,forest):
9     """
10     :param forest: forest with best parameter
11     :param X2: (n*20),testing data set in all axis (X[xi,x2])
12     :param Y2: (n*1),testing data, all data labels
13     :return: error_rate
14     """
15     s = X2.shape[0]
16     n = 0
17     Y_p = forest.predict(X2)
18     for i, y in enumerate(Y_p):
19         if (y != Y2[i]):#error_rate = #points whose predict class is different
20             n += 1
21     return n/float(s)
22
23 #read files
24 f = pd.read_csv('/Users/ccai28/Desktop/hw2_data.2.txt',sep='\t', header = None
25               , skiprows = 1)
26 #read training data
27 X = f.loc[:699,:19].values
28 Y = f.loc[:699,20].values
29
30 #read testing data
31 X_test = f.loc[700,:19].values
32 Y_test = f.loc[700,:20].values
33
34 # Range of n_estimators values to explore.
35 min_estimators = 10
36 max_estimators = 500
37 step = 10

```

```

38 number_estimators = (max_estimators - min_estimators)/step #how many trees are
    tried
39
40 print number_estimators
41
42 error_rate = []
43 n_estimators = []
44
45 #use Random Forest
46 #oob_score : bool (default=False) :Whether to use out-of-bag samples to
    estimate the generalization accuracy.
47 for i in range(min_estimators, max_estimators + step, step):
48     clf = RandomForestClassifier(oob_score=True)#warm_start=True
49     clf.set_params(n_estimators=i)
50     clf.fit(X, Y)
51
52     # Record the OOB error for each n_estimators=i setting.
53     oob_error = 1 - clf.oob_score_
54     n_estimators.append(i)
55     error_rate.append(oob_error)
56     print i, oob_error
57
58 last_oob_error = error_rate[number_estimators] #oob error of n = 500
59
60 stabilize_estimator = 0
61 #number of trees based on when the OOB error rate first stabilizes
62 for i in range(number_estimators):
63     diff_rate = abs((error_rate[i] - last_oob_error)/float(last_oob_error))
64     print i, diff_rate
65     if (diff_rate <=0.002):
66         stabilize_estimator = n_estimators[i]
67         print stabilize_estimator
68         break
69
70 print stabilize_estimator #the stablize tree number
71
72 # Generate the "OOB error rate" vs. "n_estimators" plot.
73 plt.plot(n_estimators, error_rate)
74
75 plt.xlim(min_estimators, max_estimators)
76 plt.xlabel("n_estimators")
77 plt.ylabel("OOB error rate")
78 plt.legend()
79 plt.show()
80
81 # Build a forest and compute the feature importances
82 forest = RandomForestClassifier(n_estimators = stabilize_estimator)
83 forest.fit(X,Y)
84
85 #importance and feature
86 importances = forest.feature_importances_
87 indices = np.argsort(importances)[::-1]
88
89 # Print the feature ranking
90 print("Feature ranking:")
91 for f in range(X_test.shape[1]):
92     print("%d. feature %d (%f)" % (f + 1, indices[f], importances[indices[f]]))
93
94 error_rate = Error_rate(X_test, Y_test, forest)
95 print('Random Forest: \nerror rate = {f} \n'.format(error_rate))

```

## 5 Problem 5

By running gradient boosting with deviance loss on the training data, the error rate is

$$0.193333 \quad (3)$$

The feature ranking is:

rank	index of variables	importances
1	9	(0.254537)
2	4	(0.241189)
3	14	(0.215034)
4	1	(0.034982)
5	10	(0.031280)
6	3	(0.023650)
7	7	(0.023579)
8	6	(0.020974)
9	2	(0.020797)
10	18	(0.020200)
11	5	(0.017700)
12	8	(0.015834)
13	0	(0.015096)
14	17	(0.012983)
15	13	(0.011701)
16	11	(0.010080)
17	19	(0.008941)
18	16	(0.008900)
19	15	(0.007335)
20	12	(0.005209)

Python codes of problem 5:

```

1 import pandas as pd
2 import numpy as np
3
4 from sklearn.ensemble import GradientBoostingClassifier
5
6 def Error_rate(X_test, Y_test, grid):
7     '''
8     :param grid: grid after fit by training data
9     :param X_test: (n*20), testing data set in all axis (X[xi,x2])
10    :param Y_test: (n*1), testing data, all data labels
11    :return: error_rate
12    '''
13    s = X_test.shape[0]
14    n = 0
15    Y_p = grid.predict(X_test)
16    for i in range(s):
17        if (Y_p[i] != Y_test[i]): #error_rate = #points whose predict class is
18                                #different form original one/#total points
19            n += 1
20    score = grid.score(X_test, Y_test)
21    return 1-score
22
23 #read files
24 f = pd.read_csv('/Users/ccai28/Desktop/hw2_data_2.txt', sep='\t', header = None
25                , skiprows = 1)
26
27 #read training data
28 X_train = f.loc[:699, :19].values
29 Y_train = f.loc[:699, 20].values

```

```

29 #read testing data
30 X_test = f.loc[700:,:19].values
31 Y_test = f.loc[700:,20].values
32
33 #use gradient boosting with deviance loss
34 grid = GradientBoostingClassifier(loss = 'deviance')
35 grid.fit(X_train,Y_train)
36
37 #importance and feature
38 importances = grid.feature_importances_
39 indices = np.argsort(importances)[::-1]
40
41 # Print the feature ranking
42 print("Feature ranking:")
43 for f in range(X_test.shape[1]):
44     print("%d. feature %d (%f)" % (f + 1, indices[f], importances[indices[f]]))
45
46 error_rate = Error_rate(X_test,Y_test,grid)
47 print('gradient boosting: \nerror rate = {f} \n'.format(error_rate))

```

## 6 Problem 6

With MARS, the error rate is

0.24333

(4)

Python codes of problem 6:

```

1 import pandas as pd
2 import numpy as np
3 import matplotlib.pyplot as plt
4
5 from pyearth import Earth
6
7 def Error_rate(X_test,Y_test,model):
8     """
9     :param X_test: (n*20),testing data set in all axis (X[xi,x2])
10    :param Y_test: (n*1),testing data, all data labels
11    :return: error_rate
12    """
13    s = X_test.shape[0]
14    n = 0
15    Y_p = model.predict(X_test)
16    y_median = median(Y_p)
17    for i in range(s):
18        if (Y_p[i] <= y_median):
19            y_p = -1
20        else:
21            y_p = 1
22        if (y_p != Y_test[i]):#error_rate = #points whose predict class is
different form original one/#total points
23            n += 1
24    return n/float(s)
25
26 def median(Y_test):
27     """
28     :param Y_test: (n*1),predict Y, all data labels
29     :return: median value of all predict Y
30     """
31     y_median = np.median(Y_test)
32     return y_median
33
34 #read files
35 f = pd.read_csv('/Users/ccai28/Desktop/hw2_data.2.txt',sep='\t', header = None
, skiprows = 1)

```

```

36
37 #read training data
38 X_train = f.loc[:699,:19].values
39 Y_train = f.loc[:699,20].values
40
41 #read testing data
42 X_test = f.loc[700:,:19].values
43 Y_test = f.loc[700:,20].values
44
45 #MARS from earth class
46 model = Earth() #MARS
47 model.fit(X_train,Y_train)
48
49 #error_rate
50 error_rate = Error_rate(X_test,Y_test,model)
51 print('MARS: \nerror rate = {:.f} \n'.format(error_rate))

```