



实 验 报 告

(2017 / 2018 学年 第 2 学期)

课程名称	机器学习导论
实验名称	Support vector machine
实验时间	2018 年 6 月 15 日
指导教师	王邦

姓名	游浩然	学号	U201515429
----	-----	----	------------

1 问题重述

- 给定数据集（文件 data1_Task.mat）如下图所示，参考 Demo 1 和 Demo 2，编程实现一个高斯核 SVM 进行分类。输出训练参数 C , σ 分别取 0.01, 0.03, 0.1, 0.3, 1, 3, 10, 30 时（共 64 组参数组合）的训练集上的准确率。
- 编程实现一个垃圾邮件 SVM 线性分类器，分别在训练集和测试集上计算准确率。训练数据文件：spamTrain.mat，要求导入数据时输出样本数和特征维度。测试数据文件：spamTest.mat，要求导入数据时输出样本数和特征维度。

2 Python Code for SVM

2.1 SVM

```
1  # -*- coding: utf-8 -*-
2  """
3  @author : Haoran You
4
5  """
6  import time
7  import numpy as np
8  import matplotlib.pyplot as plt
9
10 class SVM():
11     def name(self):
12         return 'svm classifier'
13
14     def svmTrain_SMO(self, X, y, C, kernel_function='linear', tol=1e-3, max_iter=5, **kargs):
15         """
16         :param X:
17         :param y:
18         :param C: punishment coefficient
19         :param kernel_function: type of kernel function; for nonlinear function, input K directly
20         :param tol: error-tolerant rate
21         :param max_iter: maximum iterations
22         :param kargs:
23         :return: model['kernelFunction']: kernel function type
24         :return: model['X']: support vector
25         :return: model['y']: label
26         :return: model['alpha']: corresponding lagrange parameters
27         :return: model['w'], model['b']: model parameters
28         """
29         start_time = time.clock()
30
31         m, n = X.shape
32         X = np.mat(X)
33         y = np.mat(y, dtype='float64')
34         y[np.where(y==0)] = -1
35         alpha = np.mat(np.zeros((m, 1)))
36         b = 0.0
37         E = np.mat(np.zeros((m, 1)))
38         iters = 0
39         eta = 0.0
40         L = 0.0
41         H = 0.0
42
43         if kernel_function == 'linear':
44             K = X*X.T
45         elif kernel_function == 'gaussian':
46             K = kargs['K_matrix']
47         else:
48             print('Kernel Error')
49             return None
50
51         print('Training ...', end='')
52         dots = 12
53         while iters < max_iter:
54             num_changed_alpha = 0
55             for i in range(m):
56                 E[i] = b + np.sum(np.multiply(np.multiply(alpha, y), K[:, i])) - y[i]
57                 if (y[i]*E[i] < -tol and alpha[i] < C) or (y[i]*E[i] > tol and alpha[i] > 0):
58                     j = np.random.randint(m)
59                     while j == i:
60                         j = np.random.randint(m)
61                 E[j] = b + np.sum(np.multiply(np.multiply(alpha, y), K[:, j])) - y[j]
62
63                 alpha_i_old = alpha[i].copy()
64                 alpha_j_old = alpha[j].copy()
```

```

66         if y[i] == y[j]:
67             L = max(0, alpha[j] + alpha[i] - C)
68             H = min(C, alpha[j] + alpha[i])
69         else:
70             L = max(0, alpha[j] - alpha[i])
71             H = min(C, C + alpha[j] - alpha[i])
72
73         if L == H:
74             continue
75
76         eta = 2 * K[i, j] - K[i, i] - K[j, j]
77         if eta >= 0:
78             continue
79
80         alpha[j] = alpha[j] - (y[j]*(E[i] - E[j])) / eta
81         alpha[j] = min(H, alpha[j])
82         alpha[j] = max(L, alpha[j])
83
84         if abs(alpha[j] - alpha_j_old) < tol:
85             alpha[j] = alpha_j_old
86             continue
87
88         alpha[i] = alpha[i] + y[i] * y[j] * (alpha_j_old - alpha[j])
89
90         b1 = b - E[i]\
91             - y[i] * (alpha[i] - alpha_i_old) * K[i, j]\
92             - y[j] * (alpha[j] - alpha_j_old) * K[i, j]
93
94         b2 = b - E[j]\
95             - y[i] * (alpha[i] - alpha_i_old) * K[i, j]\
96             - y[j] * (alpha[j] - alpha_j_old) * K[i, j]
97
98         if (0 < alpha[i] and alpha[i] < C):
99             b = b1
100         elif (0 < alpha[j] and alpha[j] < C):
101             b = b2
102         else:
103             b = (b1 + b2) / 2.0
104
105         num_changed_alpha = num_changed_alpha + 1
106
107         if num_changed_alpha == 0:
108             iters = iters + 1
109         else:
110             iters = 0
111
112         print('.', end='')
113         dots = dots + 1
114         if dots > 78:
115             dots = 0
116             print()
117
118     print('Done', end='')
119     end_time = time.clock()
120     print('(' + str(end_time - start_time) + 's)')
121     print()
122
123     idx = np.where(alpha > 0)
124     model = {
125         'X': X[idx[0], :],
126         'y': y[idx],
127         'kernelFunction': str(kernel_function),
128         'b': b,
129         'alpha': alpha[idx],
130         'w': (np.multiply(alpha, y).T * X).T
131     }
132     return model
133
134 def visualizeBoundaryLinear(self, X, y, model, title=None):
135     fig = plt.figure()
136     ax = fig.add_subplot(111)
137
138     w = model['w']
139     b = model['b']
140     xp = np.linspace(min(X[:, 0]), max(X[:, 0]), 100)
141     yp = np.squeeze(np.array(-(w[0] * xp + b) / w[1])))
142
143     ax.plot(xp, yp)
144
145     # scatter
146     X_pos = []
147     X_neg = []
148
149     sampleArray = np.concatenate((X, y), axis=1)
150     for array in list(sampleArray):
151         if array[-1]:
152             X_pos.append(array)
153         else:

```

```

154         X_neg.append(array)
155
156     X_pos = np.array(X_pos)
157     X_neg = np.array(X_neg)
158
159     if title: ax.set_title(title)
160
161     pos = plt.scatter(X_pos[:, 0], X_pos[:, 1], marker='+', c='b')
162     neg = plt.scatter(X_neg[:, 0], X_neg[:, 1], marker='o', c='y')
163
164     plt.legend((pos, neg), ('positive', 'negative'), loc=2)
165
166     plt.show()
167
168     def gaussianKernelSub(self, x1, x2, sigma):
169         X1 = np.mat(x1).reshape(-1, 1)
170         X2 = np.mat(x2).reshape(-1, 1)
171         n = -(x1-x2).T * (x1 - x2) / (2*sigma**2)
172         return np.exp(n)
173
174     def gaussianKernel(self, X, sigma):
175         start = time.clock()
176
177         print('GaussianKernel Computing ...', end='')
178         m = X.shape[0]
179         X = np.mat(X)
180         K = np.mat(np.zeros((m, m)))
181         dots = 280
182         for i in range(m):
183             if dots % 10 == 0: print('.', end='')
184             dots = dots + 1
185             if dots > 780:
186                 dots = 0
187                 print()
188             for j in range(m):
189                 K[i, j] = self.gaussianKernelSub(X[i, :].T, X[j, :].T, sigma)
190                 K[j, i] = K[i, j].copy()
191
192         print('Done', end='')
193         end = time.clock()
194         print('(' + str(end - start) + 's)')
195         print()
196         return K
197
198     def svmPredict(self, model, X, *arg):
199         m = X.shape[0]
200         p = np.mat(np.zeros((m, 1)))
201         pred = np.mat(np.zeros((m, 1)))
202
203         if model['kernelFunction'] == 'linear':
204             p = X * model['w'] + model['b']
205         else:
206             for i in range(m):
207                 prediction = 0
208                 for j in range(model['X'].shape[0]):
209                     prediction += model['alpha'][j, i] * model['y'][j, i] * \
210                         self.gaussianKernelSub(X[i, :].T, model['X'][j, :].T, *arg)
211
212                 p[i] = prediction + model['b']
213
214         pred[np.where(p >= 0)] = 1
215         pred[np.where(p < 0)] = 0
216
217         return pred
218
219     def visualizeBoundaryGaussian(self, X, y, model, sigma):
220         fig = plt.figure()
221         ax = fig.add_subplot(111)
222
223         x1plot = np.linspace(min(X[:, 0]), max(X[:, 0]), 100)
224         x2plot = np.linspace(min(X[:, 1]), max(X[:, 1]), 100)
225         X1, X2 = np.meshgrid(x1plot, x2plot)
226         X1 = np.mat(X1)
227         X2 = np.mat(X2)
228         vals = np.mat(np.zeros(X1.shape))
229
230         print('Predicting ...', end='')
231         dots = 14
232         for i in range(X1.shape[1]):
233             print('.', end='')
234             dots += 1
235             if dots == 78:
236                 dots = 0
237                 print()
238             this_X = np.concatenate((X1[:, i], X2[:, i]), axis=1)
239             vals[:, i] = self.svmPredict(model, this_X, sigma)
240         print('Done')
241

```

```

242     ax.contour(X1, X2, vals, colors='black')
243     # scatter
244     X_pos = []
245     X_neg = []
246
247     sampleArray = np.concatenate((X, y), axis=1)
248     for array in list(sampleArray):
249         if array[-1]:
250             X_pos.append(array)
251         else:
252             X_neg.append(array)
253
254     X_pos = np.array(X_pos)
255     X_neg = np.array(X_neg)
256
257     pos = plt.scatter(X_pos[:, 0], X_pos[:, 1], marker='+', c='b')
258     neg = plt.scatter(X_neg[:, 0], X_neg[:, 1], marker='o', c='y')
259
260     plt.legend((pos, neg), ('positive', 'negative'), loc=2)
261
262     plt.show()

```

2.2 data

```

1  #-*- coding: utf-8 -*-
2  """
3  @author : Haoran You
4
5  """
6  from scipy.io import loadmat
7  import numpy as np
8  import matplotlib.pyplot as plt
9
10 def loadData(filename):
11     dataDict = loadmat(filename)
12     return dataDict['X'], dataDict['y']
13
14 def plotData(X, y, title=None):
15     X_pos = []
16     X_neg = []
17
18     sampleArray = np.concatenate((X, y), axis=1)
19     for array in list(sampleArray):
20         if array[-1]:
21             X_pos.append(array)
22         else:
23             X_neg.append(array)
24
25     X_pos = np.array(X_pos)
26     X_neg = np.array(X_neg)
27
28     fig = plt.figure()
29
30     ax = fig.add_subplot(111)
31
32     if title: ax.set_title(title)
33
34     pos = plt.scatter(X_pos[:, 0], X_pos[:, 1], marker='+', c='b')
35     neg = plt.scatter(X_neg[:, 0], X_neg[:, 1], marker='o', c='y')
36
37     plt.legend((pos, neg), ('positive', 'negative'), loc=2)
38
39     plt.show()

```

2.3 Main

```

1  #-*- coding: utf-8 -*-
2  """
3  @author : Haoran You
4
5  """
6  import os
7  from data import *
8  from SVM import SVM
9
10 file_path = os.getcwd()
11
12 def Task1():
13     # load and visualize data
14     svm = SVM()
15     print('Loading and Visualizing Data ...')
16     X, y = loadData(file_path + '\data1_Task.mat')
17     plotData(X, y, title='Task 1')
18
19     print('Program paused. Press enter to continue.')
20     input()
21
22     # Training
23     print('Training SVM with RBF kernel ...')
24     param = [0.01, 0.03, 0.1, 0.3, 1, 3, 10, 30]
25     acc = []
26     for C in param:
27         for sigma in param:
28             model = svm.svmTrain_SMO(X, y, C, kernel_function='gaussian', K_matrix=svm.gaussianKernel(X, sigma))
29             pred = svm.svmPredict(model, np.mat(X), sigma)
30             acc.append(1 - np.sum(abs(pred - y)) / len(y))
31
32     print(acc)
33     f = open('results.txt', 'a+')
34     f.write('sigma\t%.3f\t%.3f\t%.3f\t%.3f\t%.3f\t%.3f\t%.3f\t%.3f\nC\n' \
35            % (0.010, 0.030, 0.100, 0.300, 1.000, 3.000, 10.000, 30.000))
36
37     count = 0
38     for i in range(len(param)):
39         content = '%.3f\t' % param[i]
40         for j in range(len(param)):
41             content += '%.3f\t' % acc[count]
42             count += 1
43         content += '\n'
44         f.write(content)
45     f.close()
46
47     # results visualizing
48     # print('Results visualizing ...')
49     # svm.visualizeBoundaryGaussian(X, y, model, sigma)
50
51 def Task2():
52     # load data
53     print('Loading data ...')
54     X_train, y_train = loadData(file_path + '\spamTrain.mat')
55     X_test, y_test = loadData(file_path + '\spamTest.mat')
56     # plotData(X_train, y_train, title='Task 2')
57
58     print('Program paused. Press enter to continue.\n')
59     input()
60
61     # training
62     svm = SVM()
63     print('Number of samples: {}'.format(X_train.shape[0]))
64     print('Number of features: {}'.format(X_train.shape[1]))
65     print('Training Linear SVM ...')
66     C = 1; sigma = 0.01;
67     model = svm.svmTrain_SMO(X_train, y_train, C, max_iter=20)
68     pred_train = svm.svmPredict(model, np.mat(X_train), sigma)
69     acc_train = 1 - abs(np.sum(pred_train - y_train)) / len(y_train)
70     print('Train accuracy: {}'.format(acc_train))
71
72     # test
73     print('Number of samples: {}'.format(X_test.shape[0]))
74     print('Number of features: {}'.format(X_test.shape[1]))
75     pred_test = svm.svmPredict(model, np.mat(X_test), sigma)
76     acc_test = 1 - abs(np.sum(pred_test - y_test)) / len(y_test)
77     print('Test accuracy: {}'.format(acc_test))
78
79 if __name__ == '__main__':
80     Task2()

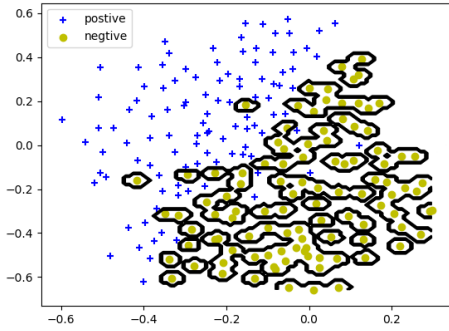
```

2.4 Results

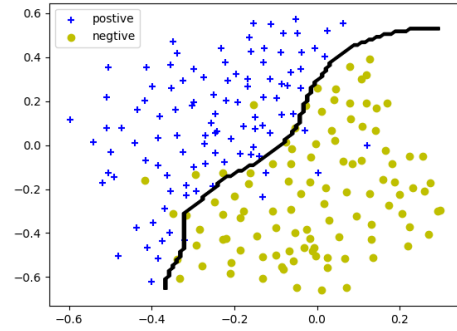
2.4.1 Task1

Table 1: Results for different parameters C & sigma

sigma C	0.01	0.03	0.1	0.3	1	3	10	30
0.01	0.498	0.498	0.502	0.502	0.502	0.498	0.502	0.502
0.03	0.502	0.502	0.502	0.867	0.502	0.498	0.834	0.498
0.1	0.498	0.502	0.943	0.872	0.834	0.81	0.815	0.81
0.3	0.502	0.981	0.948	0.905	0.867	0.81	0.498	0.502
1	1	0.995	0.948	0.934	0.905	0.848	0.815	0.829
3	1	1	0.943	0.943	0.915	0.863	0.825	0.502
10	1	1	0.962	0.934	0.929	0.9	0.848	0.815
30	1	1	0.948	0.943	0.924	0.919	0.867	0.791



C = 1 & sigma = 0.01



C = 1 & sigma = 0.1

Figure 1: Boundary Visualization

2.4.2 Task2

- Train accuracy: 0.98775
 - Number of samples: 4000
 - Number of features: 1899
- Test accuracy: 0.982
 - Number of samples: 1000
 - Number of features: 1899