

《机器学习实战》学习笔记（六）：支持向量机

原创 我是管小亮 2019-08-18 16:22:12 2755 收藏 10

版权

分类专栏: Machine Learning 文章标签: 机器学习 机器学习实战 读书笔记及代码 支持向量机

欢迎关注WX公众号：【程序员管小亮】

【机器学习】《机器学习实战》读书笔记及代码 总目录

- <https://blog.csdn.net/TeFuirnever/article/details/99701256>

GitHub代码地址：

- <https://github.com/TeFuirnever/Machine-Learning-in-Action>

目录

欢迎关注WX公众号：【程序员管小亮】

本章内容

- 1、基于最大间隔分隔数据
- 2、寻找最大间隔
- 3、简化版 SMO 算法
- 4、加速优化的完整版 Platt SMO 算法
- 5、在复杂数据上应用核函数
- 6、SVM实现手写数字识别
- 7、Sklearn构建SVM分类器
- 8、sklearn.svm.SVC
- 9、总结

关于什么是SVM? ? ?

参考文章

本章内容

- 简单介绍支持向量机
- 利用SMO进行优化
- 利用核函数对数据进行空间转换
- 将SVM和其他分类器进行对比

“由于理解支持向量机（Support Vector Machines，SVM）需要掌握一些理论知识，而这对于读者来说有一定难度，于是建议读者直接下载LIBSVM使用。”这句话我在很多书上看到相关介绍，包括机器学习实战中的理论也不是很多，西瓜书相对多一些但是不够详细，所以想要仔细看理论的话，可以移步：

- 【机器学习】《机器学习》周志华西瓜书读书笔记：第6章 - 支持向量机
- 【机器学习】《机器学习》周志华西瓜书习题参考答案：第6章 - 支持向量机

关于现在还有必要对SVM深入学习吗？这个引用一个知乎高赞回答。

窃以为，学东西，特别是科学知识不应以“火不火”作为评价标准。而是应以点带面、举一反三的建立学术知识体系，简单单纯来说也就是积累。不要看现在SVM被DL盖过了风头就产生“学了会没用，会吃亏”的念头。实际上，SVM整套体系，不论是其理论基础还是具体实现（如经典的libsvm和liblinear），都有其研究的价值。万变不离其宗，万事万物的“道”是相通的，说不准以后会在解决某些问题的时候借鉴到曾经研究过的SVM相关思想或实现技巧，彼时，便会产生融会贯通之畅快感。

闲言少叙，开始正题，本文主要关注 **序列最小优化（Sequential Minimal Optimization，SMO）** 算法，一种求解支持向量机二次规划的算法。

1、基于最大间隔分隔数据

- 优点：泛化错误率低，计算开销不大，结果易解释。
- 缺点：对参数调节和核函数的选择敏感，原始分类器不加修改仅适用于处理二类问题。
- 适用数据类型：数值型和标称型数据。

在介绍SVM这个主题之前，先解释几个概念。

- **线性可分**：可以很容易就在数据中给出一条直线将两组数据点分开。

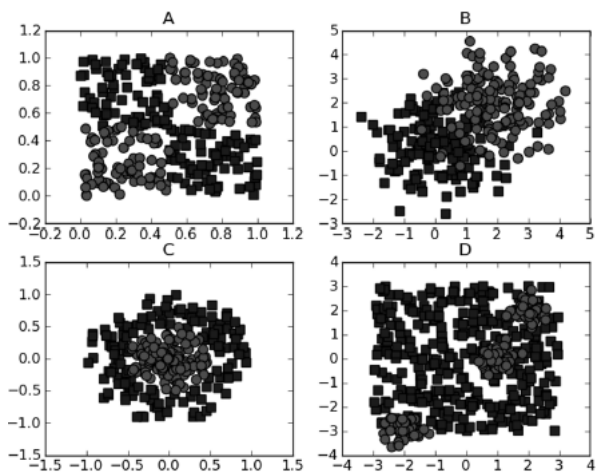


图6-1 4个线性不可分的数据集

上图的数据都是混合在一起，也就是不能用一条直线进行分类的数据，所以也就是线性不可分数据。

- **分隔超平面**：将数据集分割开来的直线。

在上面的例子中，由于数据点都在二维平面上，所以此时分隔超平面就只是一条直线。但是，如果所给的数据集是三维的，那么此时用来分隔数据的就是一个平面。显而易见，更高维的情况可以依此类推。如果数据集是1024维的，那么就需要一个1023维的某某对象来对数据进行分隔。这个1023维的某某对象到底应该叫什么？N-1维呢？该对象被称为超平面（hyperplane），也就是分类的决策边界。

理想状态是分布在超平面一侧的所有数据都属于某个类别，而分布在另一侧的所有数据则属于另一个类别。

- **间隔**：离分隔超平面最近的点，到分隔面的距离。

间隔应该尽可能地大，这是因为如果我们犯错或者在有限数据上训练分类器的话，大的间隔可以增加分类器的鲁棒性。

- **支持向量**：离分隔超平面最近的那些点。

支持向量到分割面的距离应该最大化。

2、寻找最大间隔

如何求解数据集的最佳分隔直线？

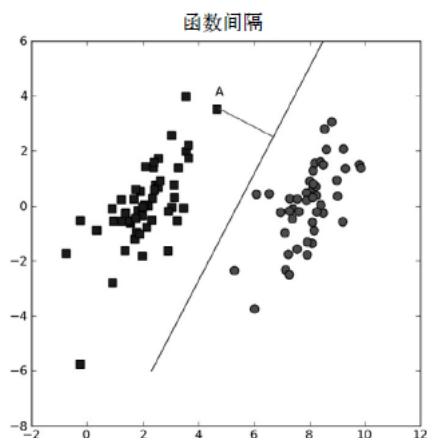


图6-3 点A到分隔平面的距离就是该点到分隔面的法线长度

推导公式为：

$$|\mathbf{w}^T \mathbf{A} + b| / \|\mathbf{w}\|$$

最大化间隔的目标就是找出分类器定义中的 w 和 b 。为此，我们必须找到具有最小间隔的数据点，而这些数据点也就是前面提到的支持向量。一旦找到具有最小间隔的数据点，我们就需要对该间隔最大化。这就可以写作：

$$\arg \max_{w,b} \left\{ \min_n (\text{label} \cdot (\mathbf{w}^T \mathbf{x} + b)) \cdot \frac{1}{\|\mathbf{w}\|} \right\}$$

直接求解上述问题相当困难，所以我们将它转换成为另一种更容易求解的形式。

$$\max_{\alpha} \left[\sum_{i=1}^m \alpha - \frac{1}{2} \sum_{i,j=1}^m \text{label}^{(i)} \cdot \text{label}^{(j)} \cdot a_i \cdot a_j \langle \mathbf{x}^{(i)}, \mathbf{x}^{(j)} \rangle \right]$$

约束条件为：

$$C \geq \alpha \geq 0, \text{ 和 } \sum_{i=1}^m \alpha_i \cdot \text{label}^{(i)} = 0$$

其中常数 C 用于控制“最大化间隔”和“保证大部分点的函数间隔小于1.0”这两个目标的权重。在优化算法的实现代码中，常数 C 是一个参数，因此可以通过调节该参数得到不同的结果。一旦求出了所有的 α ，那么分隔超平面就可以通过这些 α 来表达。

到目前为止，已经了解了一些理论知识，但是比较理论总归要回到实践上，这也是每个算法的归宿所在，通过编程，在数据集上将这些理论付诸实践。

SVM的一般流程

- (1) 收集数据：可以使用任意方法。
- (2) 准备数据：需要数值型数据。
- (3) 分析数据：有助于可视化分隔超平面。
- (4) 训练算法：SVM的大部分时间都源自训练，该过程主要实现两个参数的调优。
- (5) 测试算法：十分简单的计算过程就可以实现。
- (6) 使用算法：几乎所有分类问题都可以使用SVM，值得一提的是，SVM本身是一个二类分类器，对多类问题应用SVM需要对代码做一些修改。

3、简化版 SMO 算法

简化版SMO算法，省略了确定要优化的最佳 α 对的步骤，而是首先在数据集上进行遍历每一个 α ，再在剩余的数据集中找到另一个 α ，构成要优化的 α 对，同时对其进行优化，这里的同时是要确保公式： $\sum \alpha_i * \text{label}^{(i)} = 0$ 。

所以改变一个 α 显然会导致等式失效，所以这里需要同时改变两个 α 。接下来看实际的代码：

```

1  # -*- coding:UTF-8 -*-
2  from time import sleep
3  import matplotlib.pyplot as plt
4  import numpy as np
5  import random
6  import types
7
8  """
9  函数说明:读取数据
10
11  Parameters:
12      fileName - 文件名
13  Returns:
14      dataMat - 数据矩阵
15      labelMat - 数据标签
16  """
17  def loadDataSet(fileName):
18      dataMat = []; labelMat = []
19      fr = open(fileName)
20      for line in fr.readlines():
21          lineArr = line.strip().split('\t')
22          dataMat.append([float(lineArr[0]), float(lineArr[1])])
23          labelMat.append(float(lineArr[2]))
24      return dataMat, labelMat
25
26  """
27  函数说明:随机选择alpha
28
29  Parameters:
30      i - alpha
31      m - alpha参数个数
32  Returns:
33      j -
34  """
35  def selectJrand(i, m):
36      j = i
37      while (j == i):
38          j = int(random.uniform(0, m))
39      return j
40
41  """
42  函数说明:修剪alpha
43
44  Parameters:
45      aj - alpha值
46      H - alpha上限
47      L - alpha下限
48  Returns:
49      aj - alpah值
50  """
51  def clipAlpha(aj, H, L):
52      if aj > H:
53          aj = H
54      if L > aj:
55          aj = L
56      return aj
57
58  """
59  函数说明:简化版SMO算法
60
61  Parameters:
62      dataMatIn - 数据矩阵
63      classLabels - 数据标签
64      C - 松弛变量
65      toler - 容错率
66      maxIter - 最大迭代次数
67  Returns:
68      无
69  """
70  def smoSimple(dataMatIn, classLabels, C, toler, maxIter):
71      #转换为numpy的mat存储
72      dataMatrix = np.mat(dataMatIn); labelMat = np.mat(classLabels).transpose()
73      #初始化b参数, 统计dataMatrix的维度
74      b = 0; m, n = np.shape(dataMatrix)
75      #初始化alpha参数, 设为0
76      alphas = np.mat(np.zeros((m, 1)))
77      #初始化迭代次数

```

```

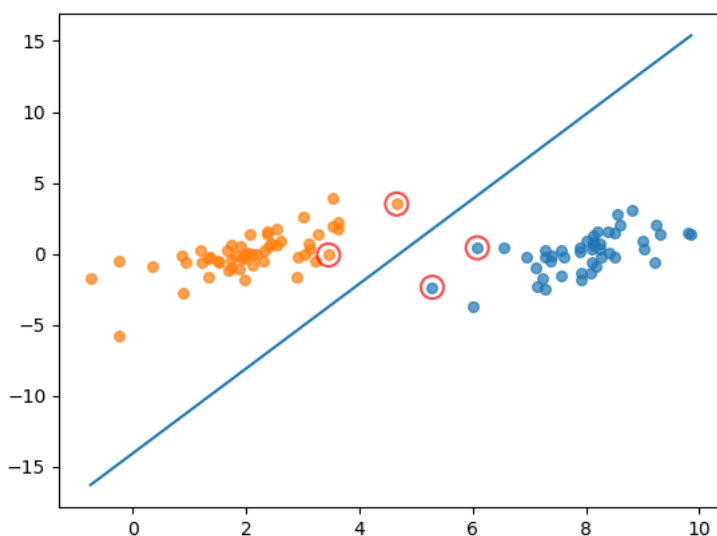
71 iter_num = 0
72 #最多迭代maxIter次
73 while (iter_num < maxIter):
74     alphaPairsChanged = 0
75     for i in range(m):
76         #步骤1: 计算误差Ei
77         fXi = float(np.multiply(alphas, labelMat).T*(dataMatrix*dataMatrix[i,:].T) + b
78         Ei = fXi - float(labelMat[i])
79         #优化alpha, 更设定一定的容错率.
80         if ((labelMat[i]*Ei < -toler) and (alphas[i] < C)) or ((labelMat[i]*Ei > toler) and (alphas[i] > 0)):
81             #随机选择另一个与alpha_i成对优化的alpha_j
82             j = selectJrand(i,m)
83             #步骤1: 计算误差Ej
84             fXj = float(np.multiply(alphas, labelMat).T*(dataMatrix*dataMatrix[j,:].T) + b
85             Ej = fXj - float(labelMat[j])
86             #保存更新前的alpha值, 使用深拷贝
87             alphaIold = alphas[i].copy(); alphaJold = alphas[j].copy();
88             #步骤2: 计算上下界L和H
89             if (labelMat[i] != labelMat[j]):
90                 L = max(0, alphas[j] - alphas[i])
91                 H = min(C, C + alphas[j] - alphas[i])
92             else:
93                 L = max(0, alphas[j] + alphas[i] - C)
94                 H = min(C, alphas[j] + alphas[i])
95             if L==H: print("L==H"); continue
96             #步骤3: 计算eta
97             eta = 2.0 * dataMatrix[i,:]*dataMatrix[j,:].T - dataMatrix[i,:]*dataMatrix[i,:].T - dataMatrix[j,:]*dataMatrix[j,:]
98             if eta >= 0: print("eta>=0"); continue
99             #步骤4: 更新alpha_j
100             alphas[j] -= labelMat[j]*(Ei - Ej)/eta
101             #步骤5: 修剪alpha_j
102             alphas[j] = clipAlpha(alphas[j],H,L)
103             if (abs(alphas[j] - alphaJold) < 0.00001): print("alpha_j变化太小"); continue
104             #步骤6: 更新alpha_i
105             alphas[i] += labelMat[j]*labelMat[i]*(alphaJold - alphas[j])
106             #步骤7: 更新b1和b2
107             b1 = b - Ei - labelMat[i]*(alphas[i]-alphaIold)*dataMatrix[i,:]*dataMatrix[i,:].T - labelMat[j]*(alphas[j]-alphaJold)*dataMatrix[j,:]*dataMatrix[i,:].T
108             b2 = b - Ej - labelMat[i]*(alphas[i]-alphaIold)*dataMatrix[i,:]*dataMatrix[j,:].T - labelMat[j]*(alphas[j]-alphaJold)*dataMatrix[j,:]*dataMatrix[j,:].T
109             #步骤8: 根据b1和b2更新b
110             if (0 < alphas[i]) and (C > alphas[i]): b = b1
111             elif (0 < alphas[j]) and (C > alphas[j]): b = b2
112             else: b = (b1 + b2)/2.0
113             #统计优化次数
114             alphaPairsChanged += 1
115             #打印统计信息
116             print("第%d次迭代 样本:%d, alpha优化次数:%d" % (iter_num,i,alphaPairsChanged))
117             #更新迭代次数
118             if (alphaPairsChanged == 0): iter_num += 1
119             else: iter_num = 0
120             print("迭代次数: %d" % iter_num)
121     return b, alphas
122
123 """
124 函数说明:分类结果可视化
125 Parameters:
126     dataMat - 数据矩阵
127     w - 直线法向量
128     b - 直线解决
129 Returns:
130     无
131 """
132 def showClassifier(dataMat, w, b):
133     #绘制样本点
134     data_plus = [] #正样本
135     data_minus = [] #负样本
136     for i in range(len(dataMat)):
137         if labelMat[i] > 0:
138             data_plus.append(dataMat[i])
139         else:
140             data_minus.append(dataMat[i])
141     data_plus_np = np.array(data_plus) #转换为numpy矩阵
142     data_minus_np = np.array(data_minus) #转换为numpy矩阵

```

```

142 plt.scatter(np.transpose(data_plus_np)[0], np.transpose(data_plus_np)[1], s=30, alpha=0.7) #正样本散点图
143 plt.scatter(np.transpose(data_minus_np)[0], np.transpose(data_minus_np)[1], s=30, alpha=0.7) #负样本散点图
144 #绘制直线
145 x1 = max(dataMat)[0]
146 x2 = min(dataMat)[0]
147 a1, a2 = w
148 b = float(b)
149 a1 = float(a1[0])
150 a2 = float(a2[0])
151 y1, y2 = (-b - a1*x1)/a2, (-b - a1*x2)/a2
152 plt.plot([x1, x2], [y1, y2])
153 #找出支持向量点
154 for i, alpha in enumerate(alphas):
155     if abs(alpha) > 0:
156         x, y = dataMat[i]
157         plt.scatter([x], [y], s=150, c='none', alpha=0.7, linewidth=1.5, edgecolor='red')
158 plt.show()
159
160 """
161 函数说明:计算w
162
163 Parameters:
164     dataMat - 数据矩阵
165     labelMat - 数据标签
166     alphas - alphas值
167 Returns:
168     无
169 """
170 def get_w(dataMat, labelMat, alphas):
171     alphas, dataMat, labelMat = np.array(alphas), np.array(dataMat), np.array(labelMat)
172     w = np.dot((np.tile(labelMat.reshape(1, -1).T, (1, 2)) * dataMat).T, alphas)
173     return w.tolist()
174
175 if __name__ == '__main__':
176     dataMat, labelMat = loadDataSet('testSet.txt')
177     b, alphas = smoSimple(dataMat, labelMat, 0.6, 0.001, 40)
178     w = get_w(dataMat, labelMat, alphas)
179     showClassifier(dataMat, w, b)
180

```



<https://blog.csdn.net/TeFuimever>

通过前面的设置进行运行程序并测试时间，我们发现大概是5s左右，虽然看起来不太差，但是不要忘了这只是一个仅有100个点的小规模数据集而已，这就意味着在更大的数据集上，收敛时间会变得更长，所以我们将通过完整的SMO算法进行加速。

4、加速优化的完整版 Platt SMO 算法

在这两个版本（简化版和完整版）中，实现alpha 的更改和代数运算的优化环节一模一样。在优化过程中，唯一的不同就是 **选择alpha的方式**。完整版的Platt SMO算法应用了一些能够提速的启发方法。

Platt SMO算法是通过一个 **外循环** 来选择第一个alpha值的，并且其选择过程会在两种方式之间进行交替：一种方式是在所有数据集上进行单遍扫描，另一种方式则是在非边界alpha中实现单遍扫描。而所谓非边界alpha指的就是那些不等于边界0或C的alpha值。对整个数据集的扫描相当容易，而实现非边界alpha值的扫描时，首先需要建立这些alpha值的列表，然后再对这个表进行遍历。同时，该步骤会跳过那些已知的不会改变的alpha值。

在选择第一个alpha值后，算法会通过一个 **内循环** 来选择第二个alpha值。在优化过程中，会通过 **最大化步长** 的方式来获得第二个alpha值。在简化版SMO算法中，我们会在选择j 之后计算错误率 E_j 。但在这里，我们会建立一个全局的缓存用于保存误差值，并从中选择使得步长或者说 $E_i - E_j$ 最大的alpha 值。

```
1  import matplotlib.pyplot as plt
2  import numpy as np
3  import random
4
5  """
6  Parameters:
7      dataMatIn - 数据矩阵
8      classLabels - 数据标签
9      C - 松弛变量
10     toler - 容错率
11 """
12 # 数据结构，维护所有需要操作的值（书上说是用于清理代码的数据结构）
13 class optStruct:
14     def __init__(self, dataMatIn, classLabels, C, toler):
15         self.X = dataMatIn#数据矩阵
16         self.labelMat = classLabels#数据标签
17         self.C = C#松弛变量
18         self.tol = toler#容错率
19         self.m = np.shape(dataMatIn)[0]#数据矩阵行数
20         self.alphas = np.mat(np.zeros((self.m,1)))#根据矩阵行数初始化alpha参数为0
21         self.b = 0#初始化b参数为0
22         #根据矩阵行数初始化误差缓存，第一列为是否有效的标志位，第二列为实际的误差E的值。
23         self.eCache = np.mat(np.zeros((self.m,2)))
24
25     """
26     Parameters:
27         fileName - 文件名
28     Returns:
29         dataMat - 数据矩阵
30         labelMat - 数据标签
31     """
32     # 读取数据
33     def loadDataSet(fileName):
34         dataMat = []; labelMat = []
35         fr = open(fileName)
36         for line in fr.readlines():#逐行读取，滤除空格等
37             lineArr = line.strip().split('\t')
38             dataMat.append([float(lineArr[0]), float(lineArr[1])])#添加数据
39             labelMat.append(float(lineArr[2]))#添加标签
40         return dataMat,labelMat
41
42     """
43     Parameters:
44         oS - 数据结构
45         k - 标号为k的数据
46     Returns:
47         Ek - 标号为k的数据误差
48     """
49     # 计算误差
50     def calcEk(oS, k):
51         fXk = float(np.multiply(oS.alphas,oS.labelMat).T*(oS.X*oS.X[k,:].T) + oS.b)
52         Ek = fXk - float(oS.labelMat[k])
53         return Ek
54
55     """
56     Parameters:
57         i - alpha_i的索引值
58         m - alpha参数个数
59     Returns:
60         j - alpha_j的索引值
61     """
62     # 函数说明：随机选择alpha_j的索引值
```

```

59 def selectJrand(i, m):
60     j = i#选择一个不等于i的
61     while (j == i):
62         j = int(random.uniform(0, m))
63     return j
64
65 """
66 Parameters:
67     i - 标号为i的数据的索引值
68     oS - 数据结构
69     Ei - 标号为i的数据误差
70 Returns:
71     j, maxK - 标号为j或maxK的数据的索引值
72     Ej - 标号为j的数据误差
73 """
74 # 内循环启发方式2
75 def selectJ(i, oS, Ei):
76     maxK = -1; maxDeltaE = 0; Ej = 0#初始化
77     oS.eCache[i] = [1,Ei]#根据i更新误差缓存
78     validEcacheList = np.nonzero(oS.eCache[:,0].A)[0]#返回误差不为0的数据的索引值
79     if (len(validEcacheList)) > 1:#有不为0的误差
80         for k in validEcacheList:#遍历,找到最大的Ek
81             if k == i: continue#不计算i,浪费时间
82             Ek = calcEk(oS, k)#计算Ek
83             deltaE = abs(Ei - Ek)#计算|Ei-Ek|
84             if (deltaE > maxDeltaE):#找到maxDeltaE
85                 maxK = k; maxDeltaE = deltaE; Ej = Ek
86         return maxK, Ej#返回maxK, Ej
87     else:#没有不为0的误差
88         j = selectJrand(i, oS.m)#随机选择alpha_j的索引值
89         Ej = calcEk(oS, j)#计算Ej
90     return j, Ej#j, Ej
91
92 """
93 Parameters:
94     oS - 数据结构
95     k - 标号为k的数据的索引值
96 Returns:
97     无
98 """
99 # 计算Ek,并更新误差缓存
100 def updateEk(oS, k):
101     Ek = calcEk(oS, k)#计算Ek
102     oS.eCache[k] = [1,Ek]#更新误差缓存
103
104 """
105 Parameters:
106     aj - alpha_j的值
107     H - alpha上限
108     L - alpha下限
109 Returns:
110     aj - 修剪后的alpah_j的值
111 """
112 # 修剪alpha_j
113 def clipAlpha(aj,H,L):
114     if aj > H:
115         aj = H
116     if L > aj:
117         aj = L
118     return aj
119
120 """
121 Parameters:
122     i - 标号为i的数据的索引值
123     oS - 数据结构
124 Returns:
125     1 - 有任意一对alpha值发生变化
126     0 - 没有任意一对alpha值发生变化或变化太小
127 """
128 # 优化的SMO算法
129 def innerL(i, oS):
130     #步骤1: 计算误差Ei
131     Ei = calcEk(oS, i)
132     #优化alpha, 设定一定的容错率。
133     if ((oS.labelMat[i] * Ei < -oS.tol) and (oS.alphas[i] < oS.C)) or\
134         ((oS.labelMat[i] * Ei > oS.tol) and (oS.alphas[i] > 0)):
135         #使用内循环启发方式2 选择alpha_j, 并计算Ej

```



```

130     j, Ej = selectJ(i, oS, Ei)
131     #保存更新前的alpha值, 使用深拷贝
132     alphaJold = oS.alphas[j].copy(); alphaJold = oS.alphas[j].copy();
133     #步骤2: 计算上下界L和H
134     if (oS.labelMat[i] != oS.labelMat[j]):
135         L = max(0, oS.alphas[j] - oS.alphas[i])
136         H = min(oS.C, oS.C + oS.alphas[j] - oS.alphas[i])
137     else:
138         L = max(0, oS.alphas[j] + oS.alphas[i] - oS.C)
139         H = min(oS.C, oS.alphas[j] + oS.alphas[i])
140     if L == H:
141         print("L==H")
142         return 0
143     #步骤3: 计算eta
144     eta = 2.0 * oS.X[i,:] * oS.X[j,:].T - oS.X[i,:] * oS.X[i,:].T - oS.X[j,:] * oS.X[j,:].T
145     if eta >= 0:
146         print("eta>=0")
147         return 0
148     #步骤4: 更新alpha_j
149     oS.alphas[j] -= oS.labelMat[j] * (Ei - Ej)/eta
150     #步骤5: 修剪alpha_j
151     oS.alphas[j] = clipAlpha(oS.alphas[j], H, L)
152     #更新j至误差缓存
153     updateEk(oS, j)
154     if (abs(oS.alphas[j] - alphaJold) < 0.00001):
155         print("alpha_j变化太小")
156         return 0
157     #步骤6: 更新alpha_i
158     oS.alphas[i] += oS.labelMat[j]*oS.labelMat[i]*(alphaJold - oS.alphas[j])
159     #更新i至误差缓存
160     updateEk(oS, i)
161     #步骤7: 更新b_1和b_2
162     b1 = oS.b - Ei - oS.labelMat[i]*(oS.alphas[i]-alphaJold)*oS.X[i,:]*oS.X[i,:].T - oS.labelMat[j]*(oS.alphas[j]-alphaJold)*
163     b2 = oS.b - Ej - oS.labelMat[i]*(oS.alphas[i]-alphaJold)*oS.X[i,:]*oS.X[j,:].T - oS.labelMat[j]*(oS.alphas[j]-alphaJold)*
164     #步骤8: 根据b_1和b_2更新b
165     if (0 < oS.alphas[i]) and (oS.C > oS.alphas[i]): oS.b = b1
166     elif (0 < oS.alphas[j]) and (oS.C > oS.alphas[j]): oS.b = b2
167     else: oS.b = (b1 + b2)/2.0
168     return 1
169 else:
170     return 0
171
172 """
173 Parameters:
174     dataMatIn - 数据矩阵
175     classLabels - 数据标签
176     C - 松弛变量
177     toler - 容错率
178     maxIter - 最大迭代次数
179 Returns:
180     oS.b - SMO算法计算的b
181     oS.alphas - SMO算法计算的alphas
182 """
183 # 完整的线性SMO算法
184 def smoP(dataMatIn, classLabels, C, toler, maxIter):
185     oS = optStruct(np.mat(dataMatIn), np.mat(classLabels).transpose(), C, toler) #初始化数据结构
186     iter = 0 #初始化当前迭代次数
187     entireSet = True; alphaPairsChanged = 0
188     #遍历整个数据集都alpha也没有更新或者超过最大迭代次数, 则退出循环
189     while (iter < maxIter) and ((alphaPairsChanged > 0) or (entireSet)):
190         alphaPairsChanged = 0
191         if entireSet: #遍历整个数据集
192             for i in range(oS.m):
193                 alphaPairsChanged += innerL(i, oS) #使用优化的SMO算法
194                 print("全样本遍历: 第%d次迭代 样本: %d, alpha优化次数: %d" % (iter, i, alphaPairsChanged))
195             iter += 1
196         else: #遍历非边界值
197             nonBoundIs = np.nonzero((oS.alphas.A > 0) * (oS.alphas.A < C))[0] #遍历不在边界0和C的alpha
198             for i in nonBoundIs:
199                 alphaPairsChanged += innerL(i, oS)
200                 print("非边界遍历: 第%d次迭代 样本: %d, alpha优化次数: %d" % (iter, i, alphaPairsChanged))
201             iter += 1
202     if entireSet: #遍历一次后改为非边界遍历

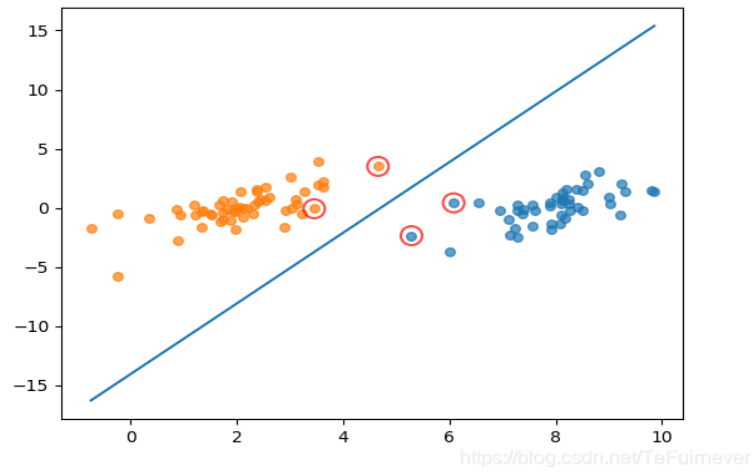
```

```

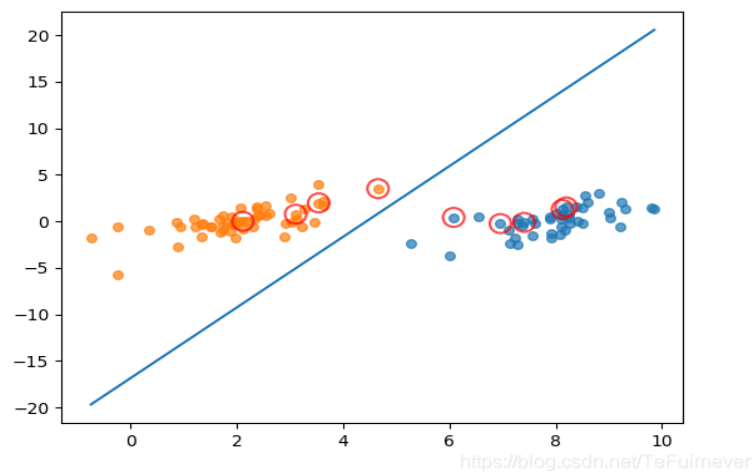
201     entireSet = False
202     elif (alphaPairsChanged == 0):#如果alpha没有更新, 计算全样本遍历
203         entireSet = True
204         print("迭代次数: %d" % iter)
205     return oS.b,oS.alphas#返回SMO算法计算的b和alphas
206
207 """
208 Parameters:
209     dataMat - 数据矩阵
210     w - 直线法向量
211     b - 直线解决
212 Returns:
213     无
214 """
215 # 分类结果可视化
216 def showClassifier(dataMat, classLabels, w, b):
217     #绘制样本点
218     data_plus = []#正样本
219     data_minus = []#负样本
220     for i in range(len(dataMat)):
221         if classLabels[i] > 0:
222             data_plus.append(dataMat[i])
223         else:
224             data_minus.append(dataMat[i])
225     data_plus_np = np.array(data_plus)#转换为numpy矩阵
226     data_minus_np = np.array(data_minus)#转换为numpy矩阵
227     plt.scatter(np.transpose(data_plus_np)[0], np.transpose(data_plus_np)[1], s=30, alpha=0.7)#正样本散点图
228     plt.scatter(np.transpose(data_minus_np)[0], np.transpose(data_minus_np)[1], s=30, alpha=0.7)#负样本散点图
229     #绘制直线
230     x1 = max(dataMat)[0]
231     x2 = min(dataMat)[0]
232     a1, a2 = w
233     b = float(b)
234     a1 = float(a1[0])
235     a2 = float(a2[0])
236     y1, y2 = (-b- a1*x1)/a2, (-b - a1*x2)/a2
237     plt.plot([x1, x2], [y1, y2])
238     #找出支持向量点
239     for i, alpha in enumerate(alphas):
240         if abs(alpha) > 0:
241             x, y = dataMat[i]
242             plt.scatter([x], [y], s=150, c='none', alpha=0.7, linewidth=1.5, edgecolor='red')
243     plt.show()
244
245 """
246 Parameters:
247     dataArr - 数据矩阵
248     classLabels - 数据标签
249     alphas - alphas值
250 Returns:
251     w - 计算得到的w
252 """
253 # 计算w
254 def calcWs(alphas,dataArr,classLabels):
255     X = np.mat(dataArr); labelMat = np.mat(classLabels).transpose()
256     m,n = np.shape(X)
257     w = np.zeros((n,1))
258     for i in range(m):
259         w += np.multiply(alphas[i]*labelMat[i],X[i,:].T)
260     return w
261
262 if __name__ == '__main__':
263     dataArr, classLabels = loadDataSet('testSet.txt')
264     b, alphas = smoP(dataArr, classLabels, 0.6, 0.001, 40)
265     w = calcWs(alphas,dataArr, classLabels)
266     showClassifier(dataArr, classLabels, w, b)
267

```

优化前:



优化后:

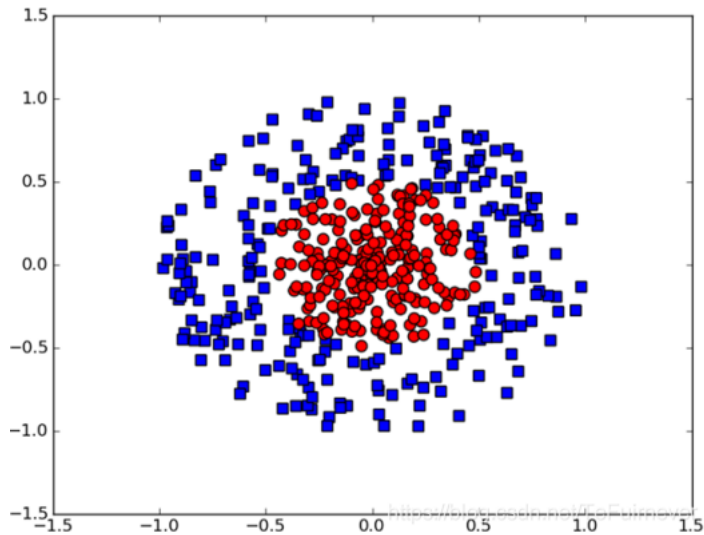


通过下面这个小程序进行计时，发现优化之后的时间相较于之前快了很多，优化前5s左右，优化后2s左右，2倍。

```
1 import time
2
3 start = time.clock()
4
5 elapsed1 = (time.clock() - start)
6 print("Time used:{:.3f}s".format(elapsed1))
```

5、在复杂数据上应用核函数

核函数的目的主要是为了解决非线性分类问题，通过核技巧将低维的非线性特征转化为高维的线性特征，从而可以通过线性模型来解决非线性的分类问题。



在图中，数据点处于一个圆中，人类的大脑能够意识到这一点。然而，对于分类器而言，它只能识别分类器的结果是大于0还是小于0。如果只在x和y轴构成的坐标系中插入直线进行分类的话，我们并不会得到理想的结果。但是或许可以对圆中的数据进行某种形式的转换，从而得到某些新的变量来表示数据。在这种表示情况下，我们就更容易得到大于0或者小于0的测试结果。在通常情况下，这种映射是通过 **核函数** 来实现的，会将低维特征空间映射到高维空间。

我们可以把核函数想象成一个 **包装器** (wrapper) 或者是 **接口** (interface)，它能把数据从某个很难处理的形式转换成为另一个较容易处理的形式。如果上述特征空间映射的说法听起来很让人迷糊的话，那么可以将它想象成为另外一种距离计算的方法。距离计算的方法有很多种，核函数一样具有多种类型。经过空间转换之后，我们可以在高维空间中解决线性问题，这也就等价于在低维空间中解决非线性问题。

接下来，我们将使用 testSetRBF.txt 和 testSetRBF2.txt，前者作为训练集，后者作为测试集。

```

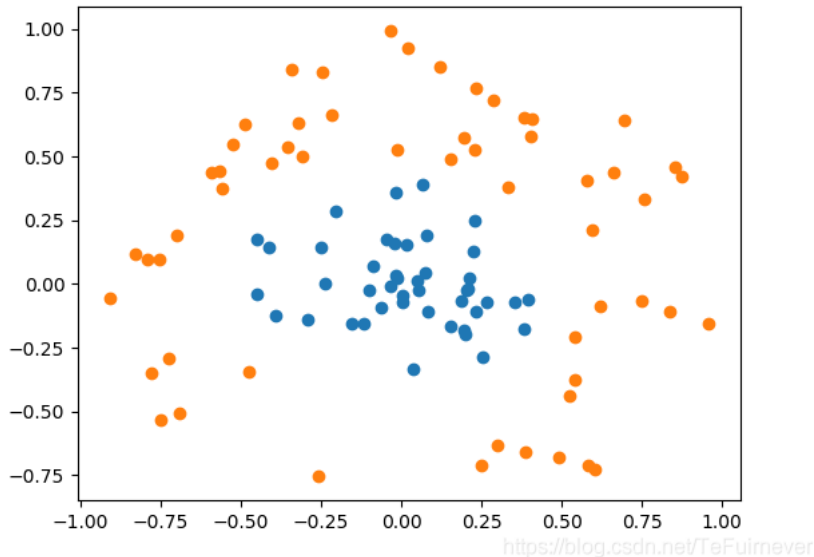
1  # -*-coding:utf-8 -*-
2  import matplotlib.pyplot as plt
3  import numpy as np
4
5  """
6  Parameters:
7      fileName - 文件名
8  Returns:
9      dataMat - 数据矩阵
10     labelMat - 数据标签
11 """
12 # 读取数据
13 def loadDataSet(fileName):
14     dataMat = []; labelMat = []
15     fr = open(fileName)
16     for line in fr.readlines():#逐行读取，滤除空格等
17         lineArr = line.strip().split('\t')
18         dataMat.append([float(lineArr[0]), float(lineArr[1])])#添加数据
19         labelMat.append(float(lineArr[2]))#添加标签
20     return dataMat,labelMat
21
22 """
23 数据可视化
24 Parameters:
25     dataMat - 数据矩阵
26     labelMat - 数据标签
27 Returns:
28     无
29 """
30 def showDataSet(dataMat, labelMat):
31     data_plus = []#正样本
32     data_minus = []#负样本
33     for i in range(len(dataMat)):
34         if labelMat[i] > 0:
35             data_plus.append(dataMat[i])
36         else:
37             data_minus.append(dataMat[i])
38     data_plus_np = np.array(data_plus)#转换为numpy矩阵
39     data_minus_np = np.array(data_minus)#转换为numpy矩阵
40     plt.scatter(np.transpose(data_plus_np)[0], np.transpose(data_plus_np)[1])#正样本散点图
41     plt.scatter(np.transpose(data_minus_np)[0], np.transpose(data_minus_np)[1])#负样本散点图

```

```

38 plt.scatter(np.transpose(data_minus_np)[0], np.transpose(data_minus_np)[1])#散点图
39 plt.show()
40
41 if __name__ == '__main__':
42     dataArr, labelArr = loadDataSet('testSetRBF.txt')#加载训练集
43     showDataSet(dataArr, labelArr)

```



可见，数据明显是线性不可分的。下面我们根据公式，编写核函数，并增加初始化参数kTup用于存储核函数有关的信息，同时我们只要将之前的内积运算变成核函数的运算即可。最后编写testRbf()函数，用于测试。创建svmMLiA.py文件，编写代码如下：

```

1  # -*-coding:utf-8 -*-
2  import matplotlib.pyplot as plt
3  import numpy as np
4  import random
5
6  """
7  Parameters:
8      dataMatIn - 数据矩阵
9      classLabels - 数据标签
10     C - 松弛变量
11     toler - 容错率
12     kTup - 包含核函数信息的元组,第一个参数存放核函数类别,第二个参数存放必要的核函数需要用到的参数
13 """
14 # 数据结构,维护所有需要操作的值
15 class optStruct:
16     def __init__(self, dataMatIn, classLabels, C, toler, kTup):
17         self.X = dataMatIn#数据矩阵
18         self.labelMat = classLabels#数据标签
19         self.C = C#松弛变量
20         self.tol = toler#容错率
21         self.m = np.shape(dataMatIn)[0]#数据矩阵行数
22         self.alphas = np.mat(np.zeros((self.m,1)))#根据矩阵行数初始化alpha参数为0
23         self.b = 0#初始化b参数为0
24         #根据矩阵行数初始化误差缓存,第一列为是否有效的标志位,第二列为实际的误差e的值。
25         self.eCache = np.mat(np.zeros((self.m,2)))
26         self.K = np.mat(np.zeros((self.m,self.m)))#初始化核
27         for i in range(self.m):#计算所有数据的核
28             self.K[:,i] = kernelTrans(self.X, self.X[i:], kTup)
29
30     """
31     Parameters:
32         X - 数据矩阵
33         A - 单个数据的向量
34         kTup - 包含核函数信息的元组
35     Returns:
36         K - 计算的核k
37     """
38
39 # 通过核函数将数据转换更高维的空间
40 def kernelTrans(X, A, kTup):

```

```

37     m, n = np.shape(X)
38     K = np.mat(np.zeros((m, 1)))
39     if kTup[0] == 'lin': K = X * A.T # 线性核函数, 只进行内积。
40     elif kTup[0] == 'rbf': # 高斯核函数, 根据高斯核函数公式进行计算
41         for j in range(m):
42             deltaRow = X[j, :] - A
43             K[j] = deltaRow * deltaRow.T
44         K = np.exp(K / (-1 * kTup[1] ** 2)) # 计算高斯核
45     else: raise NameError('核函数无法识别')
46     return K # 返回计算的核
47
48 """
49 Parameters:
50     fileName - 文件名
51 Returns:
52     dataMat - 数据矩阵
53     labelMat - 数据标签
54 """
55 # 读取数据
56 def loadDataSet(fileName):
57     dataMat = []; labelMat = []
58     fr = open(fileName)
59     for line in fr.readlines(): # 逐行读取, 滤除空格等
60         lineArr = line.strip().split('\t')
61         dataMat.append([float(lineArr[0]), float(lineArr[1])]) # 添加数据
62         labelMat.append(float(lineArr[2])) # 添加标签
63     return dataMat, labelMat
64
65 """
66 Parameters:
67     oS - 数据结构
68     k - 标号为k的数据
69 Returns:
70     Ek - 标号为k的数据误差
71 """
72 # 计算误差
73 def calcEk(oS, k):
74     fXk = float(np.multiply(oS.alphas, oS.labelMat).T * oS.K[:, k] + oS.b)
75     Ek = fXk - float(oS.labelMat[k])
76     return Ek
77
78 """
79 Parameters:
80     i - alpha_i 的索引值
81     m - alpha 参数个数
82 Returns:
83     j - alpha_j 的索引值
84 """
85 # 函数说明: 随机选择 alpha_j 的索引值
86 def selectJrand(i, m):
87     j = i # 选择一个不等于i的j
88     while (j == i):
89         j = int(random.uniform(0, m))
90     return j
91
92 """
93 Parameters:
94     i - 标号为i的数据的索引值
95     oS - 数据结构
96     Ei - 标号为i的数据误差
97 Returns:
98     j, maxK - 标号为j或maxK的数据的索引值
99     Ej - 标号为j的数据误差
100 """
101 # 内循环启发方式2
102 def selectJ(i, oS, Ei):
103     maxK = -1; maxDeltaE = 0; Ej = 0 # 初始化
104     oS.eCache[i] = [1, Ei] # 根据i更新误差缓存
105     validEcacheList = np.nonzero(oS.eCache[:, 0].A)[0] # 返回误差不为0的数据的索引值
106     if (len(validEcacheList)) > 1: # 有不等于0的误差
107         for k in validEcacheList: # 遍历, 找到最大的Ek
108             if k == i: continue # 不计算i, 浪费时间
109             Ek = calcEk(oS, k) # 计算Ek
110             deltaE = abs(Ei - Ek) # 计算|Ei - Ek|
111             if (deltaE > maxDeltaE): # 找到maxDeltaE
112                 maxK = k; maxDeltaE = deltaE; Ej = Ek
113         return maxK, Ej # 返回maxK, Ej
114     # 没有不为0的误差

```

```

108     else: # 没有个方向的误差
109         j = selectJrand(i, oS.m) # 随机选择alpha_j的索引值
110         Ej = calcEk(oS, j) # 计算Ej
111     return j, Ej#j, Ej
112
113 """
114 Parameters:
115     oS - 数据结构
116     k - 标号为k的数据的索引值
117 Returns:
118     无
119 """
120 # 计算Ek, 并更新误差缓存
121 def updateEk(oS, k):
122     Ek = calcEk(oS, k) # 计算Ek
123     oS.eCache[k] = [1, Ek] # 更新误差缓存
124
125 """
126 Parameters:
127     aj - alpha_j的值
128     H - alpha上限
129     L - alpha下限
130 Returns:
131     aj - 修剪后的alphah_j的值
132 """
133 # 修剪alpha_j
134 def clipAlpha(aj, H, L):
135     if aj > H:
136         aj = H
137     if L > aj:
138         aj = L
139     return aj
140
141 """
142 Parameters:
143     i - 标号为i的数据的索引值
144     oS - 数据结构
145 Returns:
146     1 - 有任意一对alpha值发生变化
147     0 - 没有任意一对alpha值发生变化或变化太小
148 """
149 # 优化的SMO算法
150 def innerL(i, oS):
151     # 步骤1: 计算误差i
152     Ei = calcEk(oS, i)
153     # 优化alpha, 设定一定的容错率。
154     if ((oS.labelMat[i] * Ei < -oS.tol) and (oS.alphas[i] < oS.C)) or \
155         ((oS.labelMat[i] * Ei > oS.tol) and (oS.alphas[i] > 0)):
156         # 使用内循环启发方式2选择alpha_j, 并计算Ej
157         j, Ej = selectJ(i, oS, Ei)
158         # 保存更新前的alpha值, 使用深拷贝
159         alphaIold = oS.alphas[i].copy(); alphaJold = oS.alphas[j].copy();
160         # 步骤2: 计算上下界L和H
161         if (oS.labelMat[i] != oS.labelMat[j]):
162             L = max(0, oS.alphas[j] - oS.alphas[i])
163             H = min(oS.C, oS.C + oS.alphas[j] - oS.alphas[i])
164         else:
165             L = max(0, oS.alphas[j] + oS.alphas[i] - oS.C)
166             H = min(oS.C, oS.alphas[j] + oS.alphas[i])
167         if L == H:
168             print("L==H")
169             return 0
170         # 步骤3: 计算eta
171         eta = 2.0 * oS.K[i, j] - oS.K[i, i] - oS.K[j, j]
172         if eta >= 0:
173             print("eta>=0")
174             return 0
175         # 步骤4: 更新alpha_j
176         oS.alphas[j] -= oS.labelMat[j] * (Ei - Ej)/eta
177         # 步骤5: 修剪alpha_j
178         oS.alphas[j] = clipAlpha(oS.alphas[j], H, L)
179         # 更新j至误差缓存
180         updateEk(oS, j)
181         if (abs(oS.alphas[j] - alphaJold) < 0.00001):
182             print("alpha_j变化太小")

```

```

179         return 0
180     #步骤6: 更新alpha_i
181     oS.alphas[i] += oS.labelMat[j]*oS.labelMat[i]*(alphaJold - oS.alphas[j])
182     #更新i至误差缓存
183     updateEk(oS, i)
184     #步骤7: 更新b_1和b_2
185     b1 = oS.b - Ei- oS.labelMat[i]*(oS.alphas[i]-alphaIold)*oS.K[i,i] - oS.labelMat[j]*(oS.alphas[j]-alphaJold)*oS.K[i,j]
186     b2 = oS.b - Ej- oS.labelMat[i]*(oS.alphas[i]-alphaIold)*oS.K[i,j]- oS.labelMat[j]*(oS.alphas[j]-alphaJold)*oS.K[j,j]
187     #步骤8: 根据b_1和b_2更新
188     if (0 < oS.alphas[i]) and (oS.C > oS.alphas[j]): oS.b = b1
189     elif (0 < oS.alphas[j]) and (oS.C > oS.alphas[j]): oS.b = b2
190     else: oS.b = (b1 + b2)/2.0
191     return 1
192 else:
193     return 0
194
195 """
196 Parameters:
197     dataMatIn - 数据矩阵
198     classLabels - 数据标签
199     C - 松弛变量
200     toler - 容错率
201     maxIter - 最大迭代次数
202     kTup - 包含核函数信息的元组
203 Returns:
204     oS.b - SMO算法计算的b
205     oS.alphas - SMO算法计算的alphas
206 """
207
208 # 完整的线性SMO算法
209 def smoP(dataMatIn, classLabels, C, toler, maxIter, kTup = ('lin',0)):
210     oS = optStruct(np.mat(dataMatIn), np.mat(classLabels).transpose(), C, toler, kTup)#初始化数据结构
211     iter = 0#初始化当前迭代次数
212     entireSet = True; alphaPairsChanged = 0
213     #遍历整个数据集都alpha也没有更新或者超过最大迭代次数,则退出循环
214     while (iter < maxIter) and ((alphaPairsChanged > 0) or (entireSet)):
215         alphaPairsChanged = 0
216         if entireSet:#遍历整个数据集
217             for i in range(oS.m):
218                 alphaPairsChanged += innerL(i,oS)#使用优化的SMO算法
219                 print("全样本遍历:第%d次迭代 样本:%d, alpha优化次数:%d" % (iter,i,alphaPairsChanged))
220             iter += 1
221         else:#遍历非边界值
222             nonBoundIs = np.nonzero((oS.alphas.A > 0) * (oS.alphas.A < C))[0]#遍历不在边界0和C的alpha
223             for i in nonBoundIs:
224                 alphaPairsChanged += innerL(i,oS)
225                 print("非边界遍历:第%d次迭代 样本:%d, alpha优化次数:%d" % (iter,i,alphaPairsChanged))
226             iter += 1
227         if entireSet:#遍历一次后改为非边界遍历
228             entireSet = False
229         elif (alphaPairsChanged == 0):#如果alpha没有更新,计算全样本遍历
230             entireSet = True
231         print("迭代次数: %d" % iter)
232     return oS.b,oS.alphas#返回SMO算法计算的b和alphas
233
234 """
235 Parameters:
236     k1 - 使用高斯核函数的时候表示到达率
237 Returns:
238     无
239 """
240
241 # 测试函数
242 def testRbf(k1 = 0.3):
243     dataArr,labelArr = loadDataSet('testSetRBF.txt')#加载训练集
244     b,alphas = smoP(dataArr, labelArr, 200, 0.0001, 100, ('rbf', k1))#根据训练集计算b和alphas
245     datMat = np.mat(dataArr); labelMat = np.mat(labelArr).transpose()
246     svInd = np.nonzero(alphas.A > 0)[0]#获得支持向量
247     sVs = datMat[svInd]
248     labelSV = labelMat[svInd];
249     print("支持向量个数:%d" % np.shape(sVs)[0])
250     m,n = np.shape(datMat)
251     errorCount = 0
252     for i in range(m):
253         kernelEval = kernelTrans(sVs,datMat[i,:],('rbf', k1))#计算各个点的核
254         predict = kernelEval.T * np.multiply(labelSV,alphas[svInd]) + b#根据支持向量的点, 计算超平面, 返回预测结果
255         if predict != labelArr[i]:
256             errorCount += 1
257     print("错误个数: %d" % errorCount)

```



```

250         #返回数组中各元素的正负符号, 用1和-1表示, 并统计错误个数
251         if np.sign(predict) != np.sign(labelArr[i]): errorCount += 1
252     print("训练集错误率: %.2f%%" % ((float(errorCount)/m)*100))#打印错误率
253     dataArr, labelArr = loadDataSet('testSetRBF2.txt')#加载测试集
254     errorCount = 0
255     datMat = np.mat(dataArr); labelMat = np.mat(labelArr).transpose()
256     m,n = np.shape(datMat)
257     for i in range(m):
258         kernelEval = kernelTrans(sVs, datMat[i,:], ('rbf', k1))#计算各个点的核
259         predict=kernelEval.T * np.multiply(labelSV, alphas[svInd]) + b#根据支持向量的点, 计算超平面, 返回预测结果
260         #返回数组中各元素的正负符号, 用1和-1表示, 并统计错误个数
261         if np.sign(predict) != np.sign(labelArr[i]): errorCount += 1
262     print("测试集错误率: %.2f%%" % ((float(errorCount)/m)*100))#打印错误率
263
264
265 if __name__ == '__main__':
266     testRbf()

```

```

1  >>>
2  .....
3  .....
4  迭代次数: 3
5  支持向量个数:25
6  训练集错误率: 0.00%
7  测试集错误率: 3.00%

```

可以尝试不同的K1值, 经过实验会发现K1越大, 过拟合越严重。支持向量的数目存在一个最优值。SVM的优点在于它能对数据进行高效分类。如果支持向量太少, 就可能会得到一个很差的决策边界; 如果支持向量太多, 也就相当于每次都利用整个数据集进行分类, 这种分类方法称为k近邻。

6、SVM实现手写数字识别

```

1  # -*-coding:utf-8 -*-
2  import matplotlib.pyplot as plt
3  import numpy as np
4  import random
5
6  """
7  数据结构, 维护所有需要操作的值
8  Parameters:
9      dataMatIn - 数据矩阵
10     classLabels - 数据标签
11     C - 松弛变量
12     toler - 容错率
13     kTup - 包含核函数信息的元组, 第一个参数存放核函数类别, 第二个参数存放必要的核函数需要用到的参数
14 """
15 class optStruct:
16     def __init__(self, dataMatIn, classLabels, C, toler, kTup):
17         self.X = dataMatIn#数据矩阵
18         self.labelMat = classLabels#数据标签
19         self.C = C#松弛变量
20         self.tol = toler#容错率
21         self.m = np.shape(dataMatIn)[0]#数据矩阵行数
22         self.alphas = np.mat(np.zeros((self.m,1)))#根据矩阵行数初始化alpha参数为0
23         self.b = 0#初始化b参数为0
24         #根据矩阵行数初始化误差缓存, 第一列为是否有效的标志位, 第二列为实际的误差e的值。
25         self.eCache = np.mat(np.zeros((self.m,2)))
26         self.K = np.mat(np.zeros((self.m,self.m)))#初始化核K
27         for i in range(self.m):#计算所有数据的核
28             self.K[:,i] = kernelTrans(self.X, self.X[i:], kTup)
29
30     """
31     通过核函数将数据转换更高维的空间
32     Parameters:
33         X - 数据矩阵
34         A - 单个数据的向量
35         kTup - 包含核函数信息的元组
36     Returns:
37         K - 计算的核K
38     """
39     def kernelTrans(X, A, kTup):

```

```

37     m,n = np.shape(x)
38     K = np.mat(np.zeros((m,1)))
39     if kTup[0] == 'lin': K = X * A.T#线性核函数,只进行内积。
40     elif kTup[0] == 'rbf':#高斯核函数,根据高斯核函数公式进行计算
41         for j in range(m):
42             deltaRow = X[j,:] - A
43             K[j] = deltaRow*deltaRow.T
44         K = np.exp(K/(-1*kTup[1]**2))#计算高斯核
45     else: raise NameError('核函数无法识别')
46     return K#返回计算的核
47
48 """
49 读取数据
50 Parameters:
51     fileName - 文件名
52 Returns:
53     dataMat - 数据矩阵
54     labelMat - 数据标签
55 """
56 def loadDataSet(fileName):
57     dataMat = []; labelMat = []
58     fr = open(fileName)
59     for line in fr.readlines():#逐行读取,滤除空格等
60         lineArr = line.strip().split('\t')
61         dataMat.append([float(lineArr[0]), float(lineArr[1])])#添加数据
62         labelMat.append(float(lineArr[2]))#添加标签
63     return dataMat,labelMat
64
65 """
66 计算误差
67 Parameters:
68     oS - 数据结构
69     k - 标号为k的数据
70 Returns:
71     Ek - 标号为k的数据误差
72 """
73 def calcEk(oS, k):
74     fXk = float(np.multiply(oS.alphas,oS.labelMat).T*oS.K[:,k] + oS.b)
75     Ek = fXk - float(oS.labelMat[k])
76     return Ek
77
78 """
79 函数说明:随机选择alpha_j的索引值
80 Parameters:
81     i - alpha_i的索引值
82     m - alpha参数个数
83 Returns:
84     j - alpha_j的索引值
85 """
86 def selectJrand(i, m):
87     j = i#选择一个不等于i的
88     while (j == i):
89         j = int(random.uniform(0, m))
90     return j
91
92 """
93 内循环启发方式2
94 Parameters:
95     i - 标号为i的数据的索引值
96     oS - 数据结构
97     Ei - 标号为i的数据误差
98 Returns:
99     j, maxK - 标号为j或maxK的数据的索引值
100     Ej - 标号为j的数据误差
101 """
102 def selectJ(i, oS, Ei):
103     maxK = -1; maxDeltaE = 0; Ej = 0#初始化
104     oS.eCache[i] = [1,Ei]#根据i更新误差缓存
105     validEcacheList = np.nonzero(oS.eCache[:,0].A)[0]#返回误差不为0的数据的索引值
106     if (len(validEcacheList)) > 1:#有不等于0的误差
107         for k in validEcacheList:#遍历,找到最大的Ek
108             if k == i: continue#不计算i,浪费时间
109             Ek = calcEk(oS, k)#计算Ek
110             deltaE = abs(Ei - Ek)#计算|Ei-Ek|
111             if (deltaE > maxDeltaE):#找到maxDeltaE
112                 maxK = k; maxDeltaE = deltaE; Ej = Ek
113         return maxK, Ej#返回maxK, Ej
114     #没有不为0的误差

```

```

108     else: #没有个方向的误差
109         j = selectJrand(i, oS.m) #随机选择alpha_j的索引值
110         Ej = calcEk(oS, j) #计算Ej
111     return j, Ej*j, Ej
112
113 """
114 计算Ek,并更新误差缓存
115 Parameters:
116     oS - 数据结构
117     k - 标号为k的数据的索引值
118 Returns:
119     无
120 """
121 def updateEk(oS, k):
122     Ek = calcEk(oS, k) #计算Ek
123     oS.eCache[k] = [1, Ek] #更新误差缓存
124
125 """
126 修剪alpha_j
127 Parameters:
128     aj - alpha_j的值
129     H - alpha上限
130     L - alpha下限
131 Returns:
132     aj - 修剪后的alphaj的值
133 """
134 def clipAlpha(aj, H, L):
135     if aj > H:
136         aj = H
137     if L > aj:
138         aj = L
139     return aj
140
141 """
142 优化的SMO算法
143 Parameters:
144     i - 标号为i的数据的索引值
145     oS - 数据结构
146 Returns:
147     1 - 有任意一对alpha值发生变化
148     0 - 没有任意一对alpha值发生变化或变化太小
149 """
150 def innerL(i, oS):
151     #步骤1: 计算误差Ei
152     Ei = calcEk(oS, i)
153     #优化alpha, 设定一定的容错率。
154     if ((oS.labelMat[i] * Ei < -oS.tol) and (oS.alphas[i] < oS.C)) or ((oS.labelMat[i] * Ei > oS.tol) and (oS.alphas[i] > 0)):
155         #使用内循环启发方式2选择alpha_j,并计算Ej
156         j, Ej = selectJ(i, oS, Ei)
157         #保存更新前的alpha值, 使用深拷贝
158         alphaIold = oS.alphas[i].copy(); alphaJold = oS.alphas[j].copy();
159         #步骤2: 计算上下界L和H
160         if (oS.labelMat[i] != oS.labelMat[j]):
161             L = max(0, oS.alphas[j] - oS.alphas[i])
162             H = min(oS.C, oS.C + oS.alphas[j] - oS.alphas[i])
163         else:
164             L = max(0, oS.alphas[j] + oS.alphas[i] - oS.C)
165             H = min(oS.C, oS.alphas[j] + oS.alphas[i])
166         if L == H:
167             print("L==H")
168             return 0
169         #步骤3: 计算eta
170         eta = 2.0 * oS.K[i, j] - oS.K[i, i] - oS.K[j, j]
171         if eta >= 0:
172             print("eta>=0")
173             return 0
174         #步骤4: 更新alpha_j
175         oS.alphas[j] -= oS.labelMat[j] * (Ei - Ej)/eta
176         #步骤5: 修剪alpha_j
177         oS.alphas[j] = clipAlpha(oS.alphas[j], H, L)
178         #更新j至误差缓存
179         updateEk(oS, j)
180         if (abs(oS.alphas[j] - alphaJold) < 0.00001):
181             print("alpha_j变化太小")
182             return 0
183         #步骤6: 更新alpha_i

```

```

179         oS.alphas[i] += oS.labelMat[j]*oS.labelMat[i]*(alphaJold - oS.alphas[j])
180         #更新i至误差缓存
181         updateEk(oS, i)
182         #步骤7: 更新b_1和b_2
183         b1 = oS.b - Ei- oS.labelMat[i]*(oS.alphas[i]-alphaIold)*oS.K[i,i] - oS.labelMat[j]*(oS.alphas[j]-alphaJold)*oS.K[i,j]
184         b2 = oS.b - Ej- oS.labelMat[i]*(oS.alphas[i]-alphaIold)*oS.K[i,j]- oS.labelMat[j]*(oS.alphas[j]-alphaJold)*oS.K[j,j]
185         #步骤8: 根据b_1和b_2更新b
186         if (0 < oS.alphas[i]) and (oS.C > oS.alphas[i]): oS.b = b1
187         elif (0 < oS.alphas[j]) and (oS.C > oS.alphas[j]): oS.b = b2
188         else: oS.b = (b1 + b2)/2.0
189         return 1
190     else:
191         return 0
192
193 """
194 完整的线性SMO算法
195 Parameters:
196     dataMatIn - 数据矩阵
197     classLabels - 数据标签
198     C - 松弛变量
199     toler - 容错率
200     maxIter - 最大迭代次数
201     kTup - 包含核函数信息的元组
202 Returns:
203     oS.b - SMO算法计算的b
204     oS.alphas - SMO算法计算的alphas
205 """
206
207 def smoP(dataMatIn, classLabels, C, toler, maxIter, kTup = ('lin',0)):
208     oS = optStruct(np.mat(dataMatIn), np.mat(classLabels).transpose(), C, toler, kTup)#初始化数据结构
209     iter = 0#初始化当前迭代次数
210     entireSet = True; alphaPairsChanged = 0
211     #遍历整个数据集都alpha也没有更新或者超过最大迭代次数,则退出循环
212     while (iter < maxIter) and ((alphaPairsChanged > 0) or (entireSet)):
213         alphaPairsChanged = 0
214         if entireSet:#遍历整个数据集
215             for i in range(oS.m):
216                 alphaPairsChanged += innerL(i,oS)#使用优化的SMO算法
217                 print("全样本遍历:第%d次迭代 样本:%d, alpha优化次数:%d" % (iter,i,alphaPairsChanged))
218             iter += 1
219         else:#遍历非边界值
220             nonBoundIs = np.nonzero((oS.alphas.A > 0) * (oS.alphas.A < C))[0]#遍历不在边界0和C的alpha
221             for i in nonBoundIs:
222                 alphaPairsChanged += innerL(i,oS)
223                 print("非边界遍历:第%d次迭代 样本:%d, alpha优化次数:%d" % (iter,i,alphaPairsChanged))
224             iter += 1
225         if entireSet:#遍历一次后改为非边界遍历
226             entireSet = False
227         elif (alphaPairsChanged == 0):#如果alpha没有更新,计算全样本遍历
228             entireSet = True
229         print("迭代次数: %d" % iter)
230     return oS.b,oS.alphas#返回SMO算法计算的b和alphas
231
232 """
233 将32x32的二进制图像转换为1x1024向量。
234 Parameters:
235     filename - 文件名
236 Returns:
237     returnVect - 返回的二进制图像的1x1024向量
238 """
239
240 def img2vector(filename):
241     returnVect = np.zeros((1,1024))
242     fr = open(filename)
243     for i in range(32):
244         lineStr = fr.readline()
245         for j in range(32):
246             returnVect[0,32*i+j] = int(lineStr[j])
247     return returnVect
248
249 """
250 加载图片
251 Parameters:
252     dirName - 文件夹的名字
253 Returns:
254     trainingMat - 数据矩阵
255     hwLabels - 数据标签
256 """

```

```

250 def loadImages(dirName):
251     from os import listdir
252     hwLabels = []
253     trainingFileList = listdir(dirName)
254     m = len(trainingFileList)
255     trainingMat = np.zeros((m,1024))
256     for i in range(m):
257         fileNameStr = trainingFileList[i]
258         fileStr = fileNameStr.split('.')[0]
259         classNumStr = int(fileStr.split('_')[0])
260         if classNumStr == 9: hwLabels.append(-1)
261         else: hwLabels.append(1)
262         trainingMat[i,:] = img2vector('%s/%s' % (dirName, fileNameStr))
263     return trainingMat, hwLabels
264
265 """
266 Parameters:
267     kTup - 包含核函数信息的元组
268 Returns:
269     无
270 """
271 # 测试函数
272 def testDigits(kTup=('rbf', 10)):
273     dataArr,labelArr = loadImages('trainingDigits')
274     b,alphas = smoP(dataArr, labelArr, 200, 0.0001, 10, kTup)
275     datMat = np.mat(dataArr); labelMat = np.mat(labelArr).transpose()
276     svInd = np.nonzero(alphas.A>0)[0]
277     sVs=datMat[svInd]
278     labelSV = labelMat[svInd];
279     print("支持向量个数:%d" % np.shape(sVs)[0])
280     m,n = np.shape(datMat)
281     errorCount = 0
282     for i in range(m):
283         kernelEval = kernelTrans(sVs,datMat[i,:],kTup)
284         predict=kernelEval.T * np.multiply(labelSV,alphas[svInd]) + b
285         if np.sign(predict) != np.sign(labelArr[i]): errorCount += 1
286     print("训练集错误率: %.2f%%" % (float(errorCount)/m))
287     dataArr,labelArr = loadImages('testDigits')
288     errorCount = 0
289     datMat = np.mat(dataArr); labelMat = np.mat(labelArr).transpose()
290     m,n = np.shape(datMat)
291     for i in range(m):
292         kernelEval = kernelTrans(sVs,datMat[i,:],kTup)
293         predict=kernelEval.T * np.multiply(labelSV,alphas[svInd]) + b
294         if np.sign(predict) != np.sign(labelArr[i]): errorCount += 1
295     print("测试集错误率: %.2f%%" % (float(errorCount)/m))
296
297 if __name__ == '__main__':
298     testDigits()

```

```

1 >>>
2 .....
3 .....
4 迭代次数: 7
5 支持向量个数:132
6 训练集错误率: 0.00%
7 测试集错误率: 0.01%

```

7、Sklearn构建SVM分类器

Sklearn.svm.SVC 是一个很好的模型，它是基于libsvm实现的。

```

1 # -*- coding: UTF-8 -*-
2 import numpy as np
3 import operator
4 from os import listdir
5 from sklearn.svm import SVC
6
7 """

```

```

8 Parameters:
9     filename - 文件名
10 Returns:
11     returnVect - 返回的二进制图像的1x1024向量
12 """
13 # 将32x32的二进制图像转换为1x1024向量
14 def img2vector(filename):
15     # 创建1x1024零向量
16     returnVect = np.zeros((1, 1024))
17     # 打开文件
18     fr = open(filename)
19     # 按行读取
20     for i in range(32):
21         # 读一行数据
22         lineStr = fr.readline()
23         # 每一行的前32个元素依次添加到returnVect中
24         for j in range(32):
25             returnVect[0, 32*i+j] = int(lineStr[j])
26     # 返回转换后的1x1024向量
27     return returnVect
28
29 # 手写数字分类测试
30 def handwritingClassTest():
31     # 测试集的Labels
32     hwLabels = []
33     # 返回trainingDigits目录下的文件名
34     trainingFileList = listdir('trainingDigits')
35     # 返回文件夹下文件的个数
36     m = len(trainingFileList)
37     # 初始化训练的Mat矩阵, 测试集
38     trainingMat = np.zeros((m, 1024))
39     # 从文件名中解析出训练集类别
40     for i in range(m):
41         # 获得文件的名字
42         fileNameStr = trainingFileList[i]
43         # 获得分类的数字
44         classNumber = int(fileNameStr.split('_')[0])
45         # 将获得的类别添加到hwLabels中
46         hwLabels.append(classNumber)
47         # 将每一个文件的1x1024数据存储在trainingMat矩阵中
48         trainingMat[i,:] = img2vector('trainingDigits/%s' % (fileNameStr))
49     clf = SVC(C=200, kernel='rbf')
50     clf.fit(trainingMat, hwLabels)
51     # 返回testDigits目录下的文件列表
52     testFileList = listdir('testDigits')
53     # 错误检测计数
54     errorCount = 0.0
55     # 测试数据的数量
56     mTest = len(testFileList)
57     # 从文件中解析出测试集的类别并进行分类测试
58     for i in range(mTest):
59         # 获得文件的名字
60         fileNameStr = testFileList[i]
61         # 获得分类的数字
62         classNumber = int(fileNameStr.split('_')[0])
63         # 获得测试集的1x1024向量, 用于训练
64         vectorUnderTest = img2vector('testDigits/%s' % (fileNameStr))
65         # 获得预测结果
66         # classifierResult = classify0(vectorUnderTest, trainingMat, hwLabels, 3)
67         classifierResult = clf.predict(vectorUnderTest)
68         print("分类返回结果为%d\t真实结果为%d" % (classifierResult, classNumber))
69         if(classifierResult != classNumber):
70             errorCount += 1.0
71     print("总共错了%d个数据\n错误率为%f%%" % (errorCount, errorCount/mTest * 100))
72
73 if __name__ == '__main__':
74     handwritingClassTest()
75

```

```
1 | >>>
2 | .....
3 | .....
4 | 总共错了2个数据
5 | 错误率为1.075269%
```

代码和kNN的实现是差不多的，就是换了个分类器而已。

8、sklearn.svm.SVC

sklearn.svm.SVC是一个很好的模型，SVM算法就是通过它实现的，详细的看这个博客——[sklearn.svm.SVC\(\)函数解析](#)

9、总结

支持向量机是一种分类器。之所以称为“机”是因为它会产生一个二值决策结果，即它是一种决策“机”。支持向量机的泛化错误率较低，也就是说它具有良好的学习能力，且学到的结果具有很好的推广性。这些优点使得支持向量机十分流行，有些人认为它是监督学习中最好的定式算法。

支持向量机试图通过求解一个二次优化问题来最大化分类间隔。在过去，训练支持向量机常采用非常复杂并且低效的二次规划求解方法。John Platt引入了SMO算法，此算法可以通过每次只优化2个alpha值来加快SVM的训练速度。本章首先讨论了一个简化版本所实现的SMO优化过程，接着给出了完整的Platt SMO算法。相对于简化版而言，完整版算法不仅大大地提高了优化的速度，还使其存在一些进一步提高运行速度的空间。有关这方面的工作，一个经常被引用的参考文献就是“Improvements to Platt’s SMO Algorithm for SVM Classifier Design”。

核方法或者说核技巧会将数据（有时是非线性数据）从一个低维空间映射到一个高维空间，可以将一个在低维空间中的非线性问题转换成高维空间下的线性问题来求解。核方法不止在SVM中适用，还可以用于其他算法中。而其中的径向基函数是一个常用的度量两个向量距离的核函数。支持向量机是一个二类分类器。当用其解决多类问题时，则需要额外的方法对其进行扩展。SVM的效果也对优化参数和所用核函数中的参数敏感。

下一章将通过介绍一个称为boosting的方法来结束我们有关分类的介绍。读者不久就会看到，在boosting和SVM之间存在着许多相似之处。