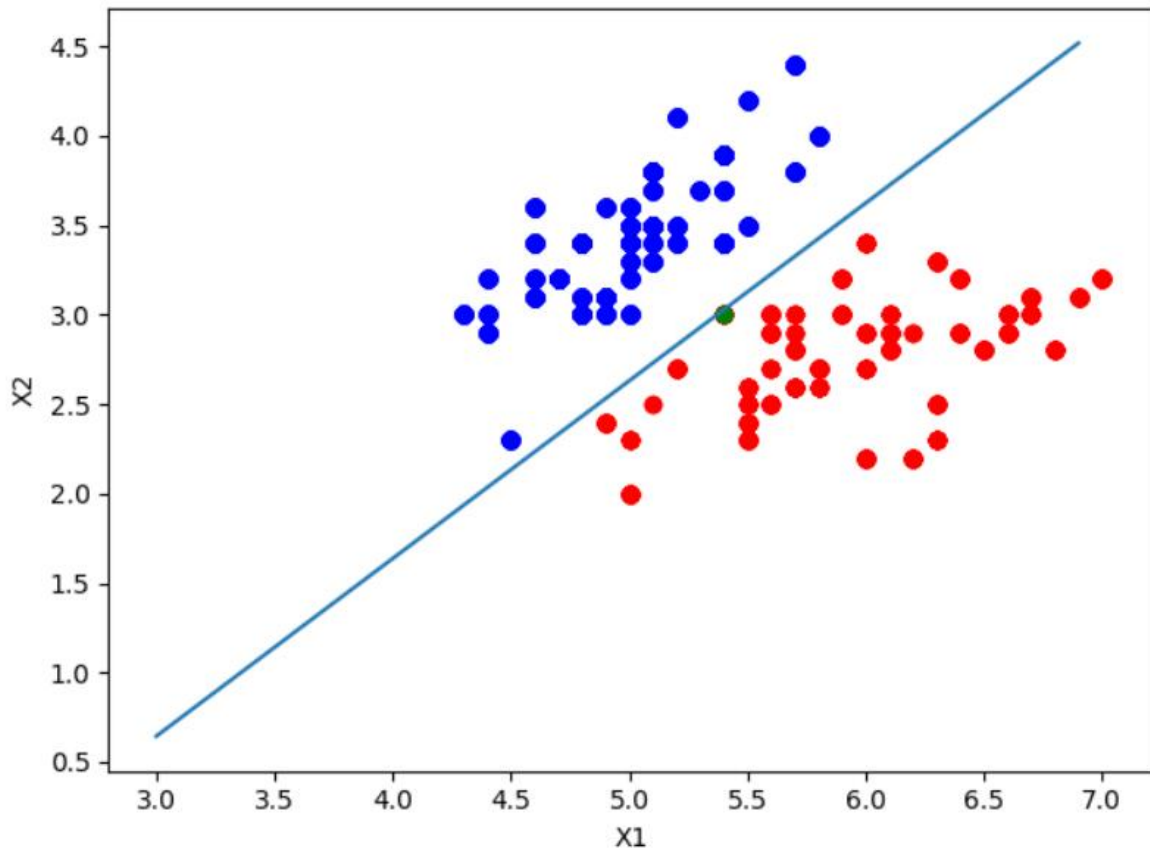# SVM鸢尾花分类问题

能达到较好的分类效果，图片如下所示



```
1  bias:      -2.481
2  w:        [ 1.05517241 -1.06206897]
3  accuracy:   1.000
```

最终可以达到100%的识别率，weight与bias如上图所示

图中红蓝两色表示两种不同的鸢尾花种类，**绿色点为支持向量**

代码中注释确为**本人手写**

`main` 函数如下

```python
1  if __name__ == "__main__":
2      trainset = create_data()
3
4      # # save data as csv when first run
5      # with open("iris_data.csv", "w") as csvdata:
6      #     writer = csv.writer(csvdata, delimiter="\n")
7      #     writer.writerows(trainset)
8
9      features, labels = trainset[:, :2], trainset[:, -1]
10
11
12     model = SVM()
13
```

```
14        # Train model
15        model.train(features, labels)
16
17        # calculate accuracy
18        y_hat = model.predict(features)
19        acc = calc_acc(labels, y_hat)
20
21        print("bias:\t\t%.3f" % (model.b))
22        print("w:\t\t" + str(model.w))
23        print("accuracy:\t%.3f" % (acc))
24
25        showpoints(trainset, model.w, model.b)
```

# 1.数据获取

使用以下python库，其中svm为自定义SVM算法文件

```
1   """
2   * @author    孟子喻
3   * @time      2021.4.30
4   * @file      classify_iris.py
5   *            svm.py
6   """
7   import numpy as np
8   import pandas as pd
9   from sklearn.datasets import load_iris
10  from svm import SVM
11  import matplotlib.pyplot as plt
```

```
1   def create_data():
2       iris = load_iris()
3       df = pd.DataFrame(iris.data, columns=iris.feature_names)
4       df['label'] = iris.target
5       df.columns = ['sepal length', 'sepal width', 'petal length', 'petal
    width', 'label']
6       data = np.array(df.iloc[:100, [0, 1, -1]])
7       for i in range(len(data)):
8           if data[i, -1] == 0:
9               data[i, -1] = -1
10      # print(data)
11      return data
12
13  # 初次使用时将数据以csv格式存储
14  # with open("iris_data.csv", "w") as csvdata:
15  #     writer = csv.writer(csvdata, delimiter="\n")
16  #     writer.writerows(trainset)
```

`create_data` 即为老师提供的数据生成函数，将其存入csv文件中供后续调用

## 2.使用SVM进行训练

我使用的svm代码结构参照了MIT的Lasse Regin Nielsen2015年所写的代码，其中有很多巧妙的设计，每一段代码都会有具体到课本公式编号的单独注释和解析

__init__

```
1   class SVM():
2       def __init__(self, max_iter=10000, kernel_type='linear', C=1.0,
    epsilon=0.001):
3           self.kernels = {
4               'linear': self.kernel_linear,
5               'quadratic': self.kernel_quadratic
6           }                                    # 根据不同的核函数实现线性与非线
    性支持向量机
7           self.max_iter = max_iter             # 最大迭代次数，超过将自动退出
8           self.kernel_type = kernel_type       # 选择核函数
9           self.C = C                           # C为惩罚参数，C越大对误分类的惩
    罚越大
10          self.epsilon = epsilon               # 设置允许差错的范围
```

把核函数 `kernel` 分为 `linear` 和 `quadratic` 两种

1. `linear`：线性核函数，对输入直接求点积
2. `quadratic`：平方核函数，对输入求点积的平方

其余参数意义见上述注释

**kernel_linear&kernel_quadratic**

```
1   # 定义核函数
2   def kernel_linear(self, x1, x2):
3       # 线性核函数
4       return np.dot(x1, x2.T)
5   def kernel_quadratic(self, x1, x2):
6       # 二次核函数
7       return (np.dot(x1, x2.T) ** 2)
```

**train解读**

```
1   # 初始化
2   n, d  = X.shape[0], X.shape[1]
3   alpha = np.zeros((n))                    # 取拉格朗日乘子初值alpha全为0
4   kernel= self.kernels[self.kernel_type]  # 设置核函数
5   count = 0                               # 计算迭代次数
```

参考课本P149，7.4.3 SMO算法，第一步取拉格朗日乘子初值alpha全为0，并设置核函数为线性（由数据点分布得）

```
1   while True:
2       count += 1
3       alpha_prev = np.copy(alpha)          # 将alpha深拷贝
4       # print(alpha.shape)
5       # print(alpha)
6
7       for j in range(0, n):
8           i = self.get_rnd_int(0, n-1, j)                    # 随机获取不同的i与
    j，得到两个优化变量
9           x_1, x_2, y_1, y_2 = X[i, :], X[j, :], y[i], y[j]   # 储存实例的特征和
    标签
10          k_ij = kernel(x_1, x_1) + kernel(x_2, x_2) - 2 * kernel(x_1, x_2)
11          if k_ij == 0:                                      # k_ij
12                  continue                                          # 保证两
    个实例不同
13                  alpha_prime_2, alpha_prime_1 = alpha[j], alpha[i]
```

随机获取两个变量的特征和标签并保证他们不同，选取alpha1和alpha2来做优化变量

```
1               (L, H) = self.compute_L_H(self.C, alpha_prime_2,
   alpha_prime_1, y_2, y_1)
2                           # print(L, H)
3                           # 计算weight和bias
4                           self.w = self.calc_w(alpha, y, X)
5                           self.b = self.calc_b(X, y, self.w)
```

这里 `compute_L_H` 是用来计算alpha2的上下边界，这个边界对应两个拉格朗日乘子alpha1和alpha2的取值范围，即[0, C] x [0, C]，其中C为惩罚系数，alpha1和alpha2即被限制在这样的一个正方形内，我们在一开始先忽略这个取值范围求最优解，求得最优情况下的取值后再与边界比较

其中求新的alpha2的上下边界的函数 `compute_L_H` 实现如下：

```
1   def compute_L_H(self, C, alpha_prime_j, alpha_prime_i, y_j, y_i):
2           # 求alpha2所在对角线端点的边界，即alpha2的取值范围
3           # print(C, alpha_prime_j, alpha_prime_i, y_j, y_i)
4           if(y_i != y_j):        # 若非同类
5               return (max(0, alpha_prime_j - alpha_prime_i), min(C, C -
   alpha_prime_i + alpha_prime_j))
6           else:                  # 若为同类
7               return (max(0, alpha_prime_i + alpha_prime_j - C), min(C,
   alpha_prime_i + alpha_prime_j))
```

由最优化问题的对 $\alpha_1$ 和 $\alpha_2$ 的约束推得，即

$\alpha_1 y_1 + \alpha_2 y_2 = $ 常数

$y_1$ 和 $y_2$ 是否相同，便对应P144 图7.8的两种情况，$\alpha_1$ 与 $\alpha_2$ 的等式约束在平行于正方形的直线上，原因如以下3公式所示，其中k也是常数，通过约束把双变量下的最优化转换为单变量下的最优化

$$y_1 \neq y_2 \text{时，} \alpha_1 - \alpha_2 = k$$

$$y_1 = y_2 \text{时，} \alpha_1 + \alpha_2 = k$$

$$L \leq \alpha_2^{new} \leq H$$

如果$y_1 \neq y_2$，则

$$L = max(0, \alpha_2^{old} - \alpha_1^{old})$$

$$H = min(C, C + \alpha_2^{old} - \alpha_1^{old})$$

如果$y_1 = y_2$，则

$$L = max(0, \alpha_2^{old} + \alpha_1^{old} - C)$$

$$H = min(C, \alpha_2^{old} + \alpha_1^{old})$$

这两个$L$与$H$的取值由以上4式决定，这里与0和C取max和min是为了裁剪保证不超出$\alpha$的取值范围

接 `train` 函数，然后是计算weight和bias

```
1                        self.w = self.calc_w(alpha, y, X)
2                        self.b = self.calc_b(X, y, self.w)
3
4                        # 计算x_i，x_j的预测值与真实值的误差
5                        E_i = self.calc_E(x_1, y_1, self.w, self.b)
6                        E_j = self.calc_E(x_2, y_2, self.w, self.b)
```

`calc_b`和`calc_w`实现如下

```
1   def calc_b(self, X, y, w):
2       b_tmp = y - np.dot(w.T, X.T)
3       return np.mean(b_tmp)
4
5   def calc_w(self, alpha, y, X):
6       return np.dot(X.T, np.multiply(alpha,y))
```

然后是计算预测值与真实值之间的误差

```
1                   E_i = self.calc_E(x_1, y_1, self.w, self.b)
2                   E_j = self.calc_E(x_2, y_2, self.w, self.b)
```

`calc_E`实现如下

```
1   def calc_E(self, x_k, y_k, w, b):
2       # 求E，即g(x)对输入x_k的预测值y_k与真实值之差
3       return self.decision_f(x_k, w, b) - y_k
4   def decision_f(self, X, w, b):
5       # 决策函数，即对输入进行预测
6       return np.sign(np.dot(w.T, X.T) + b).astype(int)
```

对应课本P145公式7.105，这里课本上写的复杂一点，但是里边一些部分已经存为weight和bias了，课本这里直接带换一下较好，`decision_f`决策函数其实就是$y = w \cdot x + b$

```
1                    # 求出alpha2未经剪辑的解
2                    alpha[j] = alpha_prime_2 + float(y_2 * (E_i - E_j))/k_ij
3                    # 利用求出的L,H对alpha2进行剪辑
4                    alpha[j] = max(alpha[j], L)
5                    alpha[j] = min(alpha[j], H)
6
7                    # 用alpha2反推alpha1
8                    alpha[i] = alpha_prime_1 + y_1*y_2 * (alpha_prime_2 -
   alpha[j])
```

这一步是正式计算 $\alpha_2^{new}$，求出取值后与前面求出的L和H边界点比较，取极大或极小值，**这里与L比较求大，与H比较求小的思路比较巧妙，能省去判断**（正常思路是如果大于H则取H，如果小于L则取L）

最后利用 $\alpha1$ 和 $\alpha2$ 的等式关系求出 $\alpha_1$，这样优化变量的整个迭代思路就很明确了，主体部分完成


完整代码

```
1    from __future__ import division, print_function
2    import os
3    import numpy as np
4    import random as rnd
5    filepath = os.path.dirname(os.path.abspath(__file__))
6
7    class SVM():
8        def __init__(self, max_iter=10000, kernel_type='linear', C=1.0,
     epsilon=0.001):
9            self.kernels = {
10               'linear': self.kernel_linear,
11               'quadratic': self.kernel_quadratic
12           }                                      # 根据不同的核函数实现线性与非线
     性支持向量机
13           self.max_iter = max_iter               # 最大迭代次数，超过将自动退出
14           self.kernel_type = kernel_type         # 选择核函数
15           self.C = C                             # C为惩罚参数，C越大对误分类的
     惩罚越大
16           self.epsilon = epsilon                 # 设置允许差错的范围
17       def train(self, X, y):
18           # 初始化
19           n, d = X.shape[0], X.shape[1]
20           alpha = np.zeros((n))                  # 取初值拉格朗日乘子alpha全为0
21           kernel = self.kernels[self.kernel_type] # 设置核函数
22           count = 0                              # 计算迭代次数
23           while True:
24               count += 1
25               alpha_prev = np.copy(alpha)        # 将alpha深拷贝
26               # print(alpha.shape)
27               # print(alpha)
28
29               for j in range(0, n):
30                   i = self.get_rnd_int(0, n-1, j)                      # 随机获
     取不同的i与j，得到两个优化变量
31                   x_1, x_2, y_1, y_2 = X[i, :], X[j, :], y[i], y[j]    # 储存实
     例的特征和标签
32                   k_ij = kernel(x_1, x_1) + kernel(x_2, x_2) - 2 *
     kernel(x_1, x_2)
```

```python
33                     if k_ij == 0:                                       # k_ij
34                         continue                                        # 保证两
个实例不同
35                     alpha_prime_2, alpha_prime_1 = alpha[j], alpha[i]
36                     # 求alpha2所在对角线端点的边界，即alpha2的取值范围
37                     (L, H) = self.compute_L_H(self.C, alpha_prime_2,
alpha_prime_1, y_2, y_1)
38                     # print(L, H)
39                     # 计算weight和bias
40                     self.w = self.calc_w(alpha, y, X)
41                     self.b = self.calc_b(X, y, self.w)
42
43                     # 计算x_i，x_j的预测值与真实值的误差
44                     E_i = self.calc_E(x_1, y_1, self.w, self.b)
45                     E_j = self.calc_E(x_2, y_2, self.w, self.b)
46
47                     # 求出alpha2未经剪辑的解
48                     alpha[j] = alpha_prime_2 + float(y_2 * (E_i - E_j))/k_ij
49                     # 利用求出的L,H对alpha2进行剪辑
50                     alpha[j] = max(alpha[j], L)
51                     alpha[j] = min(alpha[j], H)
52
53                     # 用alpha2反推alpha1
54                     alpha[i] = alpha_prime_1 + y_1*y_2 * (alpha_prime_2 -
alpha[j])
55
56                 # 检查是否超出误差允许范围
57                 diff = np.linalg.norm(alpha - alpha_prev)
58                 if diff < self.epsilon:
59                     break
60
61                 # 如果超出设定的最大迭代次数仍未求出最优解，则返回
62                 if count >= self.max_iter:
63                     print("Iteration number exceeded the max of %d iterations"
% (self.max_iter))
64                     return
65         # 临输出前计算最终的weight和bias
66         self.b = self.calc_b(X, y, self.w)
67         if self.kernel_type == 'linear':
68             self.w = self.calc_w(alpha, y, X)
69         return count
70
71     def predict(self, X):
72         return self.decision_f(X, self.w, self.b)
73     def calc_b(self, X, y, w):
74         b_tmp = y - np.dot(w.T, X.T)
75         return np.mean(b_tmp)
76
77     def calc_w(self, alpha, y, X):
78         return np.dot(X.T, np.multiply(alpha,y))
79
80     def decision_f(self, X, w, b):
81         # 决策函数，即对输入进行预测
82         return np.sign(np.dot(w.T, X.T) + b).astype(int)
83
84     def calc_E(self, x_k, y_k, w, b):
85         # 求E，即g(x)对输入x_k的预测值y_k与真实值之差
86         return self.decision_f(x_k, w, b) - y_k
```

```
 87
 88        def compute_L_H(self, C, alpha_prime_j, alpha_prime_i, y_j, y_i):
 89            # 求alpha2所在对角线端点的边界，即alpha2的取值范围
 90            # print(C, alpha_prime_j, alpha_prime_i, y_j, y_i)
 91            if(y_i != y_j):              # 若非同类
 92                return (max(0, alpha_prime_j - alpha_prime_i), min(C, C -
     alpha_prime_i + alpha_prime_j))
 93            else:                        # 若为同类
 94                return (max(0, alpha_prime_i + alpha_prime_j - C), min(C,
     alpha_prime_i + alpha_prime_j))
 95        def get_rnd_int(self, a,b,z):
 96            i = z
 97            cnt=0
 98            while i == z and cnt<1000:
 99                i = rnd.randint(a,b)
100                cnt=cnt+1
101            return i
102        # 定义核函数
103        def kernel_linear(self, x1, x2):
104            # 线性核函数
105            return np.dot(x1, x2.T)
106        def kernel_quadratic(self, x1, x2):
107            # 二次核函数
108            return (np.dot(x1, x2.T) ** 2)
```

## 3.可视化显示

最后的绘图函数也很简单了

```
 1   def showpoints(data, w, b):
 2       positive_x = []
 3       positive_y = []
 4       negative_x = []
 5       negative_y = []
 6       fig = plt.figure()
 7       ax = fig.add_subplot(111)
 8       for item in data:
 9           if item[-1] == 1:
10               positive_x.append(item[0])
11               positive_y.append(item[1])
12           if item[-1] == -1:
13               negative_x.append(item[0])
14               negative_y.append(item[1])
15           ax.scatter(positive_x, positive_y, color="r")
16           ax.scatter(negative_x, negative_y, color="b")
17       x = np.arange(3, 7, 0.1)
18       y = []
19       for item in x:
20           y.append(-w[0]/w[1]*item - b/w[1])
21       ax.plot(x, y)
22       sup_vc = get_sup_vc(data[:, 0:2], w, b)
23
24       ax.scatter(sup_vc[0], sup_vc[1], color="g")
25       plt.xlabel('X1')
26       plt.ylabel('X2')
27
28       plt.show()
```