

摘要

本篇文章参考了诸多大牛的文章写成的，对于什么是SVM做出了生动的阐述，同时也进行了线性SVM的理论推导，以及最后的编程实践，公式较多，还需静下心来一点一点



一、前言

说来惭愧，断更快半个月了，本打算是一周一篇的。感觉SVM瞬间难了不少，推导耗费了很多时间，同时身边的事情也不少，忙了许久。本篇文章参

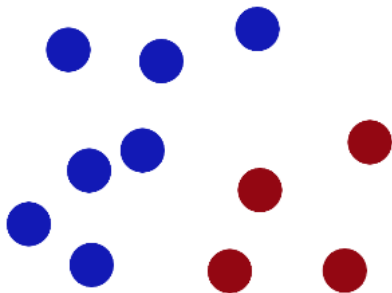
本文出现的所有代码，均可在我的github上下载，欢迎Follow、Star：<https://github.com/Jack-Cherish/Machine-Learning>

二、什么是SVM？

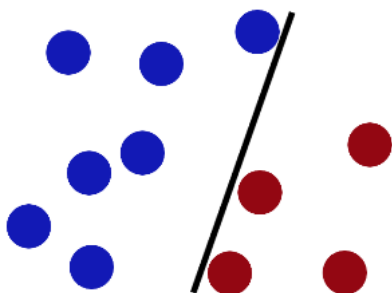
SVM的英文全称是Support Vector Machines，我们叫它支持向量机。支持向量机是我们用于分类的一种算法。让我们以一个小故事的形式，开启我

在很久以前的情人节，一位大侠要去救他的爱人，但天空中的魔鬼和他玩了一个游戏。

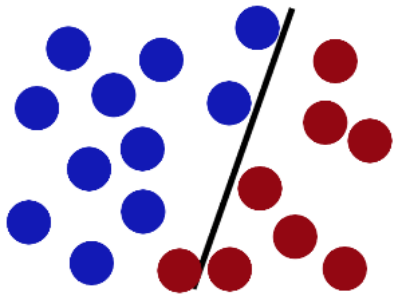
魔鬼在桌子上似乎有规律放了两种颜色的球，说："你用一根棍分开它们？要求：尽量在放更多球之后，仍然适用。"



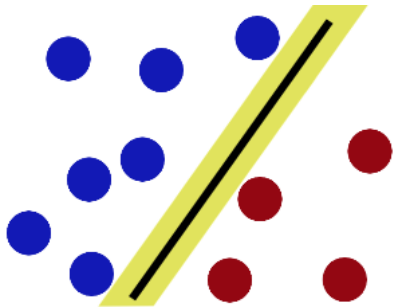
于是大侠这样放，干的不错？



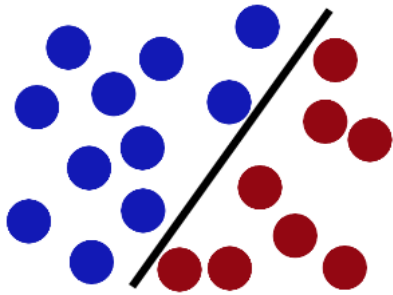
然后魔鬼，又在桌上放了更多的球，似乎有一个球站错了阵营。显然，大侠需要对棍做出调整。



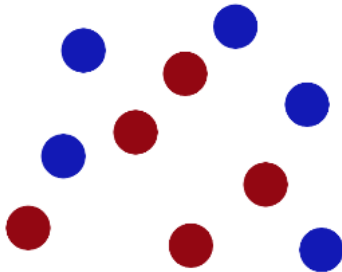
SVM就是试图把棍放在最佳位置，好让在棍的两边有尽可能大的间隙。这个间隙就是球到棍的距离。



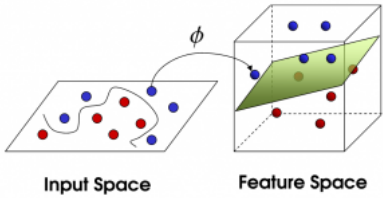
现在好了，即使魔鬼放了更多的球，棍仍然是一个好的分界线。



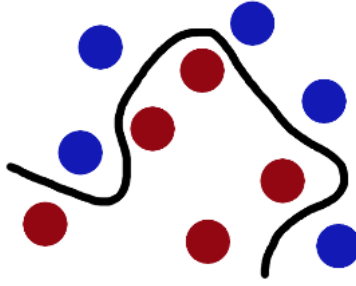
魔鬼看到大侠已经学会了一个trick(方法、招式)，于是魔鬼给了大侠一个新的挑战。



现在，大侠没有棍可以很好帮他分开两种球了，现在怎么办呢？当然像所有武侠片中一样大侠桌子一拍，球飞到空中。然后，凭借大侠的轻功，大侠到了两种球的中间。



现在，从空中的魔鬼的角度看这些球，这些球看起来像是被一条曲线分开了。



再之后，无聊的大人们，把这些球叫做**data**，把棍子叫做**classifier**，找到最大间隙的**trick**叫做**optimization**，拍桌子叫做**kernelling**，那张纸叫做**hyperplane**。更为直观地感受一下吧(需要翻墙)：<https://www.youtube.com/watch?v=3liCbRZPrZA>

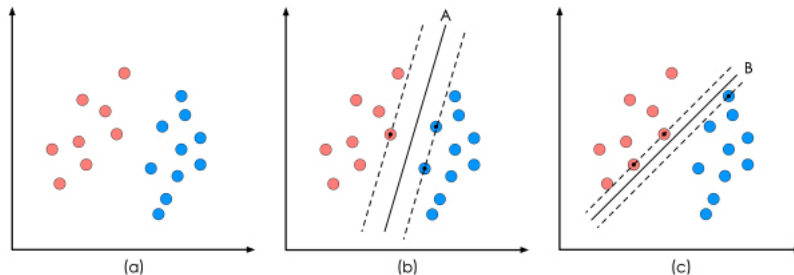
概述一下：

当一个分类问题，数据是线性可分的，也就是用一根棍就可以将两种小球分开的时候，我们只要将棍的位置放在让小球距离棍的距离最大化的位置即最大间隔的过程，就叫做最优化。但是，现实往往是很残酷的，一般的数据是线性不可分的，也就是找不到一个棍将两种小球很好的分类。这个时候，我们一样，将小球拍起，用一张纸代替小棍将小球进行分类。想要让数据飞起，我们需要的东西就是核函数(kernel)，用于切分小球的纸，就是超平面。

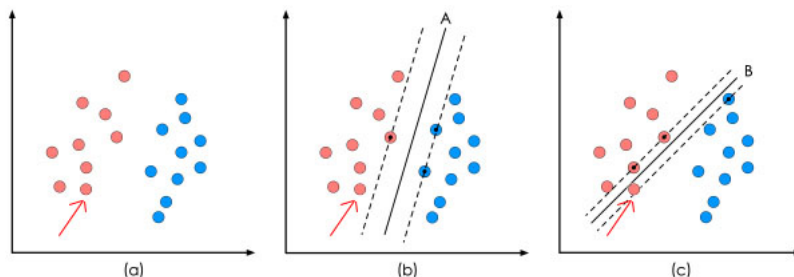
也许这个时候，你还是似懂非懂，没关系。根据刚才的描述，可以看出，问题是从线性可分延伸到线性不可分的。那么，我们就按照这个思路，进行

三、线性SVM

先看下线性可分的二分类问题。



上图中的(a)是已有的数据，红色和蓝色分别代表两个不同的类别。数据显然是线性可分的，但是将两类数据点分开的直线显然不止一条。上图的(b)和B、C两种不同的分类方案，其中黑色实线为分界线，术语称为“决策面”。每个决策面对应了一个线性分类器。虽然从分类结果上看，分类器A和分类器C都是正确的。但是他们的性能是有差距的，看下图：



在“决策面”不变的情况下，我又添加了一个红点。可以看到，分类器B依然能很好的分类结果，而分类器C则出现了分类错误。显然分类器B的“决策面”放置的位置，SVM算法也是这么认为的，它的依据就是分类器B的分类间隔比分类器C的分类间隔大。这里涉及到第一个SVM独有的“最大间隔”。在保证决策面方向不变且不会出现错分样本的情况下移动决策面，会在原来的决策面两侧找到两个极限位置（越过该位置就会产生错分现象），如虚线的位置由决策面的方向和距离原决策面最近的几个样本的位置决定。而这两条平行虚线正中间的分界线就是在保持当前决策面方向不变的前提下的最优决策面。之间的垂直距离就是这个最优决策面对应的分类间隔。显然每一个可能把数据集正确分开方向都有一个最优决策面（有些方向无论如何移动决策面的位置类样本完全分开），而不同方向的最优决策面的分类间隔通常是不同的，那个具有“最大间隔”的决策面就是SVM要寻找的最优解。而这个真正的最优解线所穿过的样本点，就是SVM中的支持样本点，称为“支持向量”。

1、数学建模

求解这个“决策面”的过程，就是最优化。一个最优化问题通常有两个基本的因素：1）目标函数，也就是你希望什么东西的什么指标达到最好；2）约束条件，也就是你希望什么东西的什么指标达到最差。在线性SVM算法中，目标函数显然就是那个“分类间隔”，而优化对象则是决策面。所以要对SVM问题进行建模，先要对上述两个对象（“分类间隔”和“决策面”）进行数学描述。按照一般的思维习惯，我们先描述决策面。

数学建模的时候，先在二维空间建模，然后再推广到多维。

(1) "决策面"方程

我们都知道二维空间下一条直线的方式如下所示：

$$y = ax + b$$

现在我们做个小小的改变，让原来的x轴变成 x_1 ，y轴变成 x_2

$$x_2 = ax_1 + b$$

移项得：

$$ax_1 - x_2 + b = 0$$

将公式向量化得：

$$\begin{bmatrix} a & -1 \end{bmatrix} \begin{bmatrix} x_1 \\ x_2 \end{bmatrix} + b = 0$$

进一步向量化，用w列向量和x列向量和标量y进一步向量化：

$$\boldsymbol{\omega}^T \boldsymbol{x} + \gamma = 0$$

其中，向量w和x分别为：

$$\boldsymbol{\omega} = [\omega_1, \omega_2]^T, \boldsymbol{x} = [x_1, x_2]^T$$

这里 $w_1=a$ ， $w_2=-1$ 。我们都知道，最初的那个直线方程a和b的几何意义，a表示直线的斜率，b表示截距，a决定了直线与x轴正方向的夹角，b决定点位置。那么向量化后的直线的w和r的几何意义是什么呢？

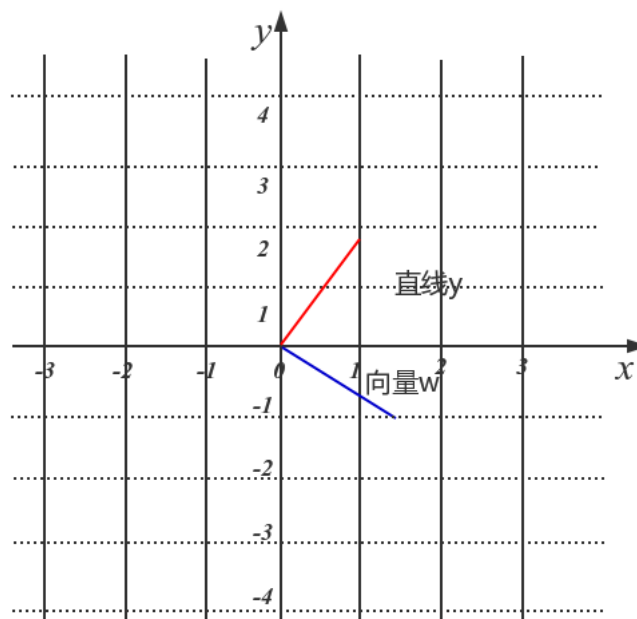
现在假设：

$$a = \sqrt{3}, b = 0$$

可得：

$$\boldsymbol{\omega} = [\sqrt{3}, -1]^T$$

在坐标轴上画出直线和向量w：



蓝色的线代表向量 w ，红色的线代表直线 y 。我们可以看到向量 w 和直线的关系为垂直关系。这说明了向量 w 也控制这直线的方向，只不过是这个直线的。标量 γ 的作用也没有变，依然决定了直线的截距。此时，我们称 w 为直线的法向量。

二维空间的直线方程已经推导完成，将其推广到 n 维空间，就变成了超平面方程。（一个超平面，在二维空间的例子就是一个直线）但是它的公式没变，

$$\omega^T x + \gamma = 0$$

不同之处在于：

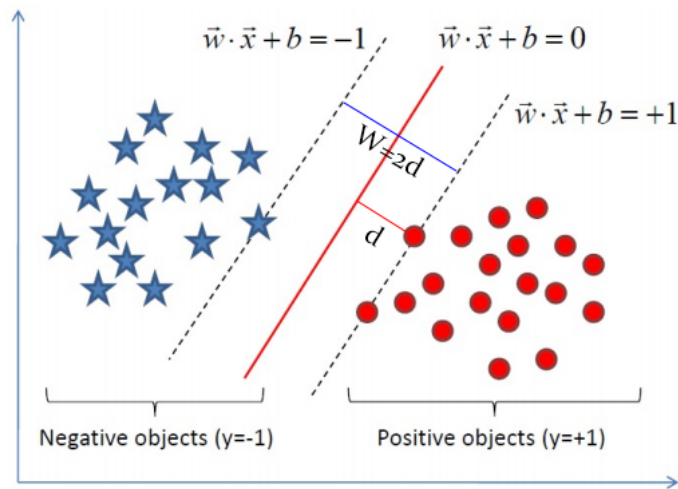
$$\omega = [\omega_1, \omega_2, \dots, \omega_n]^T$$

$$x = [x_1, x_2, \dots, x_n]^T$$

我们已经顺利推导出了“决策面”方程，它就是我们的超平面方程，之后，我们统称其为超平面方程。

（2）“分类间隔”方程

现在，我们依然对于一个二维平面的简单例子进行推导。



我们已经知道间隔的大小实际上就是支持向量对应的样本点到决策面的距离的二倍。那么图中的距离 d 我们怎么求？我们高中都学过，点到直线的距离下：

$$d = \left| \frac{Ax_0 + By_0 + C}{\sqrt{A^2 + B^2}} \right|$$

公式中的直线方程为 $Ax_0 + By_0 + C = 0$ ，点 P 的坐标为 (x_0, y_0) 。

现在，将直线方程扩展到多维，求得我们现在的超平面方程，对公式进行如下变形：

$$d = \frac{|\omega^T x + \gamma|}{\|\omega\|}$$

这个 d 就是“分类间隔”。其中 $\|\omega\|$ 表示 w 的二范数，求所有元素的平方和，然后再开方。比如对于二维平面：

$$\omega = [\omega_1, \omega_2]^T$$

那么，

$$\|\omega\| = \sqrt{\omega_1^2 + \omega_2^2}$$

我们目的是为了找出一个分类效果好的超平面作为分类器。分类器的好坏的评定依据是分类间隔 $W=2d$ 的大小，即分类间隔 w 越大，我们认为这个超平面越好。此时，求解超平面的问题就变成了求解分类间隔 W 最大化的为题。 W 的最大化也就是 d 最大化的。

（3）约束条件

看起来，我们已经顺利获得了目标函数的数学形式。但是为了解w的最大值。我们不得不面对如下问题：

- 我们如何判断超平面是否将样本点正确分类？
- 我们知道要求距离d的最大值，我们首先需要找到支持向量上的点，怎么在众多的点中选出支持向量上的点呢？

上述我们需要面对的问题就是约束条件，也就是说我们优化的变量d的取值范围受到了限制和约束。事实上约束条件一直是最优化问题里最让人头疼的。我们已经知道了这些约束条件确实存在，就不得不用数学语言对他们进行描述。但SVM算法通过一些巧妙的小技巧，将这些约束条件融合到一个不等式里

这个二维平面上有两种点，我们分别对它们进行标记：

- 红颜色的圆点标记为1，我们人为规定其为正样本；
- 蓝颜色的五角星标记为-1，我们人为规定其为负样本。

对每个样本点 x_i 加上一个类别标签 y_i ：

$$y_i = \begin{cases} +1 & \text{红色点} \\ -1 & \text{蓝色点} \end{cases}$$

如果我们的超平面方程能够完全正确地对上图的样本点进行分类，就会满足下面的方程：

$$\begin{cases} \omega^T x_i + \gamma > 0 & y_i = 1 \\ \omega^T x_i + \gamma < 0 & y_i = -1 \end{cases}$$

如果我们要求再高一点，假设决策面正好处于间隔区域的中轴线上，并且相应的支持向量对应的样本点到决策面的距离为d，那么公式进一步写成：

$$\begin{cases} \frac{\omega^T x_i + \gamma}{\|\omega\|} \geq d & \forall y_i = 1 \\ \frac{\omega^T x_i + \gamma}{\|\omega\|} \leq -d & \forall y_i = -1 \end{cases}$$

上述公式的解释就是，对于所有分类标签为1和-1样本点，它们到直线的距离都大于等于d(支持向量上的样本点到超平面的距离)。公式两边都除以d，

$$\begin{cases} \omega_d^T x_i + \gamma_d \geq 1 & \forall y_i = 1 \\ \omega_d^T x_i + \gamma_d \leq -1 & \forall y_i = -1 \end{cases}$$

其中，

$$\underline{\omega_d} = \frac{\omega}{\|\omega\|d}, \quad \underline{\gamma_d} = \frac{\gamma}{\|\omega\|d}$$

因为 $\|\omega\|$ 和d都是标量。所以上述公式的两个矢量，依然描述一条直线的法向量和截距。

$$\begin{aligned} \omega_d^T x + \gamma_d &= 0 \\ \omega^T x + \gamma &= 0 \end{aligned}$$

上述两个公式，都是描述一条直线，数学模型代表的意义是一样的。现在，让我们对 ω_d 和 γ_d 重新起个名字，就叫它们w和γ。因此，我们就可以说：'间隔的两类样本点，我们一定可以找到一些超平面，使其对于所有的样本点均满足下面的条件：'

$$\begin{cases} \omega^T x_i + \gamma \geq 1 & \forall y_i = 1 \\ \omega^T x_i + \gamma \leq -1 & \forall y_i = -1 \end{cases}$$

上述方程即给出了SVM最优化问题的约束条件。这时候，可能有人会问了，为什么标记为1和-1呢？因为这样标记方便我们将上述方程变成如下形式

$$y_i(\omega^T x_i + \gamma) \geq 1 \quad \forall x_i$$

正是因为标签为1和-1，才方便我们将约束条件变成一个约束方程，从而方便我们的计算。

（4）线性SVM优化问题基本描述

现在整合一下思路，我们已经得到我们的目标函数：

$$d = \frac{|\omega^T x + \gamma|}{\|\omega\|}$$

我们的优化目标是d最大化。我们已经说过，我们是用支持向量上的样本点求解d的最大化的问题的。那么支持向量上的样本点有什么特点呢？

$$|\omega^T x_i + \gamma| = 1 \quad \forall \text{支持向量上的样本点 } x_i$$

你赞同这个观点吗？所有支持向量上的样本点，都满足如上公式。如果不赞同，请重看“分类间隔”方程推导过程。

现在我们就可以将我们的目标函数进一步化简：

$$d = \frac{1}{\|\omega\|}$$

因为，我们只关心支持向量上的点。随后我们求解d的最大化问题变成了 $\|\omega\|$ 的最小化问题。进而 $\|\omega\|$ 的最小化问题等效于

$$\min \frac{1}{2} \|\omega\|^2$$

为什么要做这样的等效呢？这是为了在进行最优化的过程中对目标函数求导时比较方便，但这绝对不影响最优化问题最后的求解。我们将最终的目标函数和约束条件放在一起进行描述：

$$\begin{aligned} \min & \frac{1}{2} \|\omega\|^2 \\ \text{s.t. } & y_i(\omega^T x_i + b) \geq 1, i = 1, 2, \dots, n \end{aligned}$$

这里n是样本点的总个数，缩写s.t.表示"Subject to"，是“服从某某条件”的意思。上述公式描述的是一个典型的不等式约束条件下的二次型函数优化问题，即支持向量机的基本数学模型。

（5）求解准备

我们已经得到支持向量机的基本数学模型，接下来的问题就是如何根据数学模型，求得我们想要的最优解。在学习求解方法之前，我们得知道一点，求解的方法有一个前提，就是我们的目标函数必须是凸函数。理解凸函数，我们还要先明确另一个概念，凸集。在凸几何中，凸集(convex set)是在凸组空间的子集。看一幅图可能更容易理解：

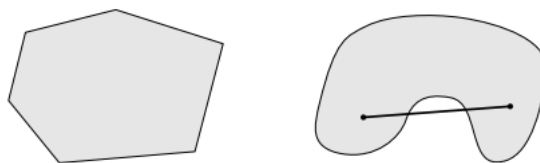


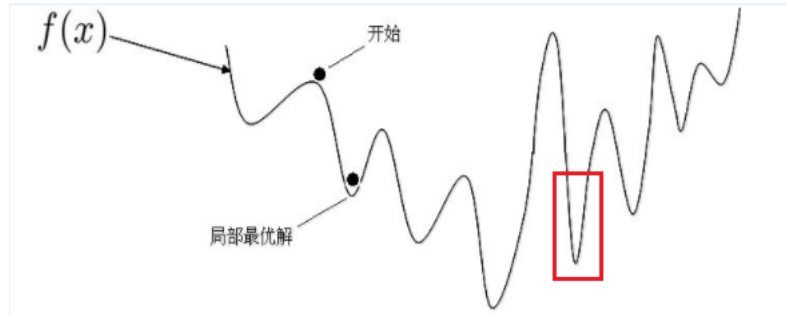
Figure 1: Examples of a convex set (a) and a non-convex set (b).

左右两图都是一个集合。如果集合中任意2个元素连线上的点也在集合中，那么这个集合就是凸集。显然，上图中的左图是一个凸集，上图中的右图是

非凸集。凸函数的定义也是如此，其几何意义表示为函数任意两点连线上的值大于对应自变量处的函数值。若这里凸集C即某个区间L，那么，设函数f为定义在L上的函数，若对L上的任意两点 x_1, x_2 和任意的实数 $\lambda, \lambda \in (0, 1)$ ，总有：

$$f(\lambda x_1 + (1 - \lambda)x_2) \leq \lambda f(x_1) + (1 - \lambda)f(x_2)$$

则函数 f 称为 L 上的凸函数，当且仅当其上图（在函数图像上方的点集）为一个凸集。再看一幅图，也许更容易理解：



像上图这样的函数，它整体就是一个非凸函数，我们无法获得全局最优解的，只能获得局部最优解。比如红框内的部分，如果单独拿出来，它就是我们的目标函数：

$$\min \frac{1}{2} \|\omega\|^2$$

很显然，它是一个凸函数。所以，可以使用我接下来讲述的方法求取最优解。

通常我们需要求解的最优化问题有如下几类：

■ 无约束优化问题，可以写为：

$$\min f(x)$$

■ 有等式约束的优化问题，可以写为：

$$\begin{aligned} \min f(x) \\ \text{s.t. } h_i(x) = 0, \quad i = 1, 2, \dots, n \end{aligned}$$

■ 有不等式约束的优化问题，可以写为：

$$\begin{aligned} \min f(x) \\ \text{s.t. } g_i(x) \leq 0, \quad i = 1, 2, \dots, n \\ h_j(x) = 0, \quad j = 1, 2, \dots, m \end{aligned}$$

对于第(a)类的优化问题，尝试使用的方法就是费马大定理(Fermat)，即使用求取函数 $f(x)$ 的导数，然后令其为零，可以求得候选最优值，再在这些候选值是凸函数，可以保证是最优解。这也就是我们高中经常使用的求函数的极值的方法。

对于第(b)类的优化问题，常常使用的方法就是拉格朗日乘子法(Lagrange Multiplier)，即把等式约束 $h_i(x)$ 用一个系数与 $f(x)$ 写为一个式子，称为拉格朗日函数。通过拉格朗日函数对各个变量求导，令其为零，可以求得候选值集合，然后验证求得最优值。

对于第(c)类的优化问题，常常使用的方法就是KKT条件。同样地，我们把所有的等式、不等式约束与 $f(x)$ 写为一个式子，也叫拉格朗日函数，系数也叫子，通过一些条件，可以求出最优值的必要条件，这个条件称为KKT条件。

必要条件和充要条件如果不理解，可以看下面这句话：

- A的必要条件就是A可以推出的结论
- A的充分条件就是可以推出A的前提

了解到这些，现在让我们再看一下我们的最优化问题：

$$\begin{aligned} \min \frac{1}{2} \|\omega\|^2 \\ \text{s.t. } y_i(\omega^T x_i + b) \geq 1, i = 1, 2, \dots, n \end{aligned}$$

现在，我们的这个对优化问题属于哪一类？很显然，它属于第(c)类问题。在学习求解最优化问题之前，我们还要学习两个东西：拉格朗日函数和KKT条件。

（6）拉格朗日函数

首先，我们先要从宏观的视野上了解一下**拉格朗日对偶问题出现的原因和背景**。

我们知道我们要求解的是最小化问题，所以一个直观的想法是如果我能够构造一个函数，使得该函数在可行解区域内与原目标函数完全一致，而在可行值非常大，甚至是无穷大，那么这个**没有约束条件的新目标函数的优化问题**就与**原来有约束条件的原始目标函数的优化问题**是等价的问题。这就是使用拉格的，它将约束条件放到目标函数中，**从而将有约束优化问题转换为无约束优化问题**。

随后，人们又发现，使用拉格朗日获得的函数，使用求导的方法求解依然困难。进而，需要对问题再进行一次转换，即使用一个数学技巧：**拉格朗日**；所以，显而易见的是，我们在拉格朗日优化我们的问题这个道路上，**需要进行下面二个步骤**：

- 将有约束的原始目标函数转换为无约束的新构造的拉格朗日目标函数
- 使用拉格朗日对偶性，将不易求解的优化问题转化为易求解的优化

下面，进行第一步：**将有约束的原始目标函数转换为无约束的新构造的拉格朗日目标函数**

公式变形如下：

$$\mathcal{L}(w, b, \alpha) = \frac{1}{2} \|w\|^2 - \sum_{i=1}^n \alpha_i (y_i (w^T x_i + b) - 1)$$

其中 α_i 是拉格朗日乘子， α_i 大于等于0，是我们构造新目标函数时引入的系数变量(我们自己设置)。现在我们令：

$$\theta(w) = \max_{\alpha_i \geq 0} \mathcal{L}(w, b, \alpha)$$

当样本点不满足约束条件时，即在**可行解区域外**：

$$y_i (w^T x_i + b) < 1$$

此时，我们将 α_i 设置为正无穷，此时 $\theta(w)$ 显然也是正无穷。

当样本点满足约束条件时，即在**可行解区域内**：

$$y_i (w^T x_i + b) \geq 1$$

此时，显然 $\theta(w)$ 为原目标函数本身。我们将上述两种情况结合一下，就得到了新的目标函数：

$$\theta(w) = \begin{cases} \frac{1}{2} \|w\|^2 & x \in \text{可区域} \\ +\infty & x \in \text{非可行区域} \end{cases}$$

此时，再看我们的初衷，就是为了建立一个在可行解区域内与原目标函数相同，在可行解区域外函数值趋近于无穷大的新函数，现在我们做到了。

现在，我们的问题变成了求新目标函数的最小值，即：

$$\min_{w, b} \theta(w) = \min_{w, b} \max_{\alpha_i \geq 0} \mathcal{L}(w, b, \alpha) = p^*$$

这里用 p^* 表示这个问题的最优值，且和最初的问题是等价的。

接下来，我们进行第二步：**将不易求解的优化问题转化为易求解的优化**

我们看一下我们的新目标函数，先求最大值，再求最小值。这样的话，我们首先就要面对带有需要求解的参数 w 和 b 的方程，而 α_i 又是不等式约束，这好做。所以，我们需要使用拉格朗日函数对偶性，将最小和最大的位置交换一下，这样就变成了：

$$\max_{\alpha_i \geq 0} \min_{w, b} \mathcal{L}(w, b, \alpha) = d^*$$

交换以后的新问题是原始问题的对偶问题，这个新问题的最优值用 d^* 来表示。而且 $d^* \leq p^*$ 。我们关心的是 $d=p$ 的时候，这才是我们要的解。需要什么 $d=p$ 呢？

- 首先必须满足这个优化问题是凸优化问题。
- 其次，需要满足KKT条件。

凸优化问题的定义是：**求取最小值的目标函数为凸函数的一类优化问题**。目标函数是凸函数我们已经知道，这个优化问题又是求最小值。所以我们的凸优化问题。

接下来，就是探讨是否满足KKT条件了。

(7) KKT条件

我们已经使用拉格朗日函数对我们的目标函数进行了处理，生成了一个新的目标函数。通过一些条件，可以求出最优值的必要条件，这个条件就是KKT条件。一个最优化模型能够表示成下列标准形式：

$$\begin{aligned} \min f(\mathbf{x}) \\ \text{s.t. } h_j(\mathbf{x}) = 0, j = 1, 2, \dots, p \\ g_k(\mathbf{x}) \leq 0, k = 1, 2, \dots, q \\ \mathbf{x} \in X \subset \mathbb{R}^n \end{aligned}$$

KKT条件的全称是Karush-Kuhn-Tucker条件，KKT条件是说最优值条件必须满足以下条件：

- 条件一：经过拉格朗日函数处理之后的新目标函数 $L(w, b, \alpha)$ 对 x 求导为零；
- 条件二： $h(x) = 0$ ；
- 条件三： $\alpha * g(x) = 0$ ；

对于我们的优化问题：

$$\begin{aligned} \min \frac{1}{2} \|\omega\|^2 \\ \text{s.t. } y_i(\omega^T x_i + b) \geq 1, i = 1, 2, \dots, n \end{aligned}$$

显然，条件二已经满足了。另外两个条件为啥也满足呢？

这里原谅我省略一系列证明步骤，感兴趣的可以移步这里：[点我查看](#)

这里已经给出了很好的解释。现在，凸优化问题和KKT都满足了，问题转换成了对偶问题。而求解这个对偶学习问题，可以分为三个步骤：首先要让 L 和 b 最小化，然后求对 α 的极大，最后利用SMO算法求解对偶问题中的拉格朗日乘子。现在，我们继续推导。

(8) 对偶问题求解

第一步：

根据上述推导已知：

$$\begin{aligned} \max_{\alpha_i \geq 0} \min_{w, b} \mathcal{L}(w, b, \alpha) = d^* \\ \mathcal{L}(w, b, \alpha) = \frac{1}{2} \|w\|^2 - \sum_{i=1}^n \alpha_i (y_i(w^T x_i + b) - 1) \end{aligned}$$

首先固定 α ，要让 $L(w, b, \alpha)$ 关于 w 和 b 最小化，我们分别对 w 和 b 偏导数，令其等于0，即：

$$\begin{aligned} \frac{\partial \mathcal{L}}{\partial w} = 0 \Rightarrow w = \sum_{i=1}^n \alpha_i y_i x_i \\ \frac{\partial \mathcal{L}}{\partial b} = 0 \Rightarrow \sum_{i=1}^n \alpha_i y_i = 0 \end{aligned}$$

将上述结果带回 $L(w, b, \alpha)$ 得到：

$$\mathcal{L}(\mathbf{w}, b, \alpha) = \frac{1}{2} \|\mathbf{w}\|^2 - \sum_{i=1}^n \alpha_i [y_i (\mathbf{w}^T \mathbf{x}_i + b) - 1]$$

从上面的最后一个式子，我们可以看出，此时的 $\mathcal{L}(\mathbf{w}, b, \alpha)$ 函数只含有一个变量，即 α_i 。

第二步：

现在内侧的最小值求解完成，我们求解外侧的最大值，从上面的式子得到

$$\begin{aligned} \max_{\alpha} \quad & \sum_{i=1}^n \alpha_i - \frac{1}{2} \sum_{i,j=1}^n \alpha_i \alpha_j y_i y_j \mathbf{x}_i^T \mathbf{x}_j \\ \text{s.t.} \quad & \alpha_i \geq 0, i = 1, 2, \dots, n \\ & \sum_{i=1}^n \alpha_i y_i = 0 \end{aligned}$$

现在我们的优化问题变成了如上的形式。对于这个问题，我们有更高效的优化算法，即序列最小优化（SMO）算法。我们通过这个优化算法能得到 α 们就可以求解出 \mathbf{w} 和 b ，进而求得我们最初的目的：找到超平面，即“决策平面”。

总结一句话：我们为啥使出吃奶的劲儿进行推导？因为我们要将最初的原始问题，转换到可以使用SMO算法求解的问题，这是一种最流行的求解方法求解方法？因为它牛逼啊！

对于上述问题，我们就可以使用SMO算法进行求解了，但是，SMO算法又是什么呢？

2、SMO算法

现在，我们已经得到了可以用SMO算法求解的目标函数，但是对于怎么编程实现SMO算法还是感觉无从下手。那么现在就聊聊如何使用SMO算法进

（1）Platt的SMO算法

1996年，John Platt发布了一个称为SMO的强大算法，用于训练SVM。SM表示序列最小化(Sequential Minimal Optimizaion)。Platt的SMO算法是分解为多个小优化问题来求解的。这些小优化问题往往很容易求解，并且对它们进行顺序求解的结果与将它们作为整体来求解的结果完全一致的。在结果收敛时，SMO算法的求解时间短很多。

SMO算法的目标是求出一系列 α 和 b ，一旦求出了这些 α ，就很容易计算出权重向量 \mathbf{w} 并得到分隔超平面。

SMO算法的工作原理是：每次循环中选择两个 α 进行优化处理。一旦找到了一对合适的 α ，那么就增大其中一个同时减小另一个。这里所谓的两个 α 必须符合以下两个条件，条件之一就是两个 α 必须要在间隔边界之外，而且第二个条件则是这两个 α 还没有进行过区间化处理或者不在其

（2）SMO算法的解法

先来定义特征到结果的输出函数为：

$$u = \omega^T x + b$$

接着，我们回忆一下原始优化问题，如下：

$$\begin{aligned} \min & \frac{1}{2} \|\omega\|^2 \\ \text{s.t.} & y_i(\omega^T x_i + b) \geq 1, i = 1, 2, \dots, n \end{aligned}$$

求导得：

$$\omega = \sum_{i=1}^n \alpha_i y_i x_i$$

将上述公式带入输出函数中：

$$u = \sum_{i=1}^n \alpha_i y_i x_i^T x + b$$

与此同时，拉格朗日对偶后得到最终的目标化函数：

$$\begin{aligned} \max_{\alpha} & \sum_{i=1}^n \alpha_i - \frac{1}{2} \sum_{i,j=1}^n \alpha_i \alpha_j y_i y_j x_i^T x_j \\ \text{s.t.} & \alpha_i \geq 0, i = 1, 2, \dots, n \\ & \sum_{i=1}^n \alpha_i y_i = 0 \end{aligned}$$

将目标函数变形，在前面增加一个符号，将最大值问题转换成最小值问题：

$$\begin{aligned} \min_{\alpha} & \frac{1}{2} \sum_{i=1}^n \sum_{j=1}^n \alpha_i \alpha_j y_i y_j x_i^T x_j - \sum_{i=1}^n \alpha_i \\ \text{s.t.} & \alpha_i \geq 0, i = 1, 2, \dots, n \\ & \sum_{i=1}^n \alpha_i y_i = 0 \end{aligned}$$

实际上，对于上述目标函数，是存在一个假设的，即数据100%线性可分。但是，目前为止，我们知道几乎所有数据都不那么“干净”。这时我们就可以的**松弛变量**(slack variable)，来允许有些数据点可以处于超平面的错误的一侧。这样我们的优化目标就能保持仍然不变，但是此时我们的约束条件有所改

$$\begin{aligned} \text{s.t.} & C \geq \alpha_i \geq 0, i = 1, 2, \dots, n \\ & \sum_{i=1}^n \alpha_i y_i = 0 \end{aligned}$$

根据KKT条件可以得出其中 α_i 取值的意义为：

$$\alpha_i = 0 \Leftrightarrow y_i u_i \geq 1$$

$$0 < \alpha_i < C \Leftrightarrow y_i u_i = 1$$

$$\alpha_i = C \Leftrightarrow y_i u_i \leq 1$$

- 对于第1种情况，表明 α_i 是正常分类，在边界内部；
- 对于第2种情况，表明 α_i 是支持向量，在边界上；
- 对于第3种情况，表明 α_i 是在两条边界之间。

而最优解需要满足KKT条件，即上述3个条件都得满足，以下几种情况出现将会不满足：

$$y_i u_i \leq 1 \quad \alpha_i < C$$

$$y_i u_i \geq 1 \quad \alpha_i > 0$$

$$y_i u_i = 1 \quad \alpha_i = 0 \text{ 或者 } \alpha_i = C$$

也就是说，如果存在不能满足KKT条件的 α_i ，那么需要更新这些 α_i ，这是第一个约束条件。此外，更新的同时还要受到第二个约束条件的限制，即：

$$\sum_{i=1}^n \alpha_i y_i = 0$$

因为这个条件，我们同时更新两个 α 值，因为只有成对更新，才能保证更新之后的值仍然满足和为0的约束，假设我们选择的两个乘子为 α_1 和 α_2 ：

$$\alpha_1^{new} y_1 + \alpha_2^{new} y_2 = \alpha_1^{old} y_1 + \alpha_2^{old} y_2 = \zeta$$

其中， ζ 为常数。因为两个因子不好同时求解，所以可以先求第二个乘子 α_2 的解（ α_2^{new} ），得到 α_2 的解（ α_2^{new} ）之后，再用 α_2 的解（ α_2^{new} ）解（ α_1^{new} ）。为了求解 α_2^{new} ，得先确定 α_2^{new} 的取值范围。假设它的上下边界分别为H和L，那么有：

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

接下来，综合下面两个条件：

$$C \geq \alpha_i \geq 0, \quad i = 1, 2, \dots, n$$

$$\alpha_1^{new} y_1 + \alpha_2^{new} y_2 = \alpha_1^{old} y_1 + \alpha_2^{old} y_2 = \zeta$$

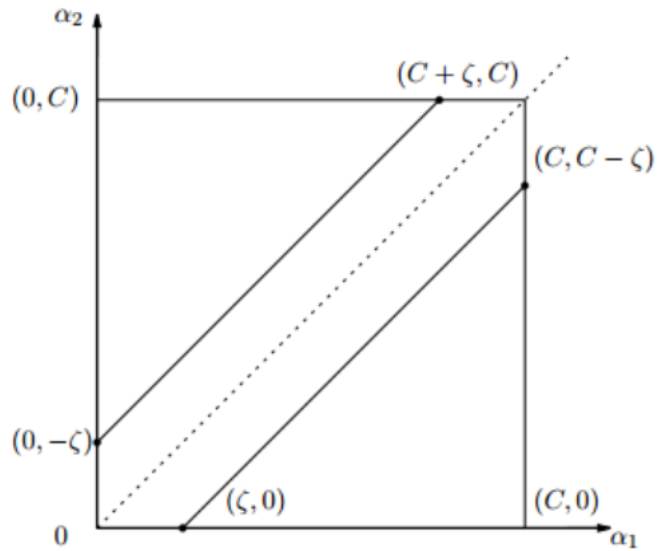
当 y_1 不等于 y_2 时，即一个为正1，一个为负1的时候，可以得到：

$$\alpha_1^{old} - \alpha_2^{old} = \zeta$$

所以有：

$$L = \max(0, -\zeta), \quad H = \min(C, C - \zeta)$$

此时，取值范围如下图所示：



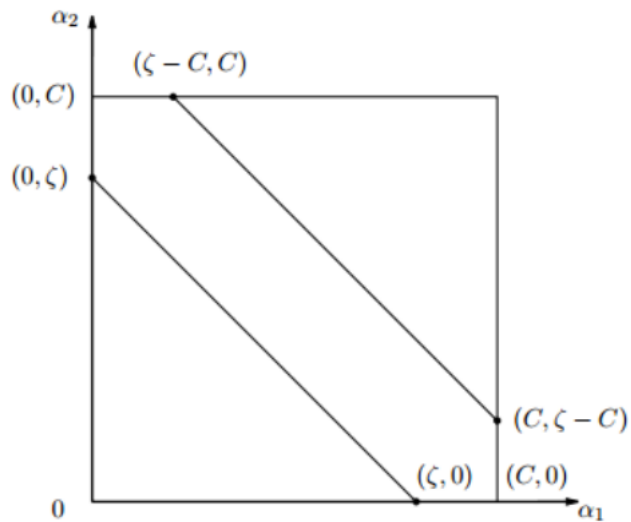
当 y_1 等于 y_2 时，即两个都为正1或者都为负1，可以得到：

$$\alpha_1^{old} + \alpha_2^{old} = \zeta$$

所以有：

$$L = \max(0, \zeta - C), \quad H = \min(C, \zeta)$$

此时，取值范围如下图所示：



如此，根据 y_1 和 y_2 异号或同号，可以得出 α_2 new的上下界分别为：

$$\begin{cases} L = \max(0, \alpha_2^{old} - \alpha_1^{old}), H = \min(C, C + \alpha_2^{old} - \alpha_1^{old}) & \text{if } y_1 \neq y_2 \\ L = \max(0, \alpha_2^{old} + \alpha_1^{old} - C), H = \min(C, \alpha_2^{old} + \alpha_1^{old}) & \text{if } y_1 = y_2 \end{cases}$$

这个界限就是编程的时候需要用到的。已经确定了边界，接下来，就是推导迭代式，用于更新 α 值。

我们已经知道，更新 α 的边界，接下来就是讨论如何更新 α 值。我们依然假设选择的两个乘子为 α_1 和 α_2 。固定这两个乘子，进行推导。于是目标函数3

$$\begin{aligned}
W(a_2) &= \sum_{i=1}^n \alpha_i - \frac{1}{2} \sum_{i=1}^n \sum_{j=1}^n y_i y_j \mathbf{x}_i^T \mathbf{x}_j \alpha_i \alpha_j \\
&= \alpha_1 + \alpha_2 + \sum_{i=3}^n \alpha_i - \frac{1}{2} \sum_{i=1}^n \left(\sum_{j=1}^2 y_i y_j \alpha_i \alpha_j \mathbf{x}_i^T \mathbf{x}_j + \sum_{j=3}^n y_i y_j \alpha_i \alpha_j \mathbf{x}_i^T \mathbf{x}_j \right) \\
&= \alpha_1 + \alpha_2 + \sum_{i=3}^n \alpha_i - \frac{1}{2} \sum_{i=1}^2 \left(\sum_{j=1}^2 y_i y_j \alpha_i \alpha_j \mathbf{x}_i^T \mathbf{x}_j + \sum_{j=3}^n y_i y_j \alpha_i \alpha_j \mathbf{x}_i^T \mathbf{x}_j \right) \\
&\quad - \frac{1}{2} \sum_{i=3}^n \left(\sum_{j=1}^2 y_i y_j \alpha_i \alpha_j \mathbf{x}_i^T \mathbf{x}_j + \sum_{j=3}^n y_i y_j \alpha_i \alpha_j \mathbf{x}_i^T \mathbf{x}_j \right) \\
&= \alpha_1 + \alpha_2 + \sum_{i=3}^n \alpha_i - \frac{1}{2} \sum_{i=1}^2 \sum_{j=1}^2 y_i y_j \alpha_i \alpha_j \mathbf{x}_i^T \mathbf{x}_j - \sum_{i=1}^2 \sum_{j=3}^n y_i y_j \alpha_i \alpha_j \mathbf{x}_i^T \mathbf{x}_j \\
&\quad - \frac{1}{2} \sum_{i=3}^n \sum_{j=3}^n y_i y_j \alpha_i \alpha_j \mathbf{x}_i^T \mathbf{x}_j \\
&= \alpha_1 + \alpha_2 - \frac{1}{2} \alpha_1^2 \mathbf{x}_1^T \mathbf{x}_1 - \frac{1}{2} \alpha_2^2 \mathbf{x}_2^T \mathbf{x}_2 - y_1 y_2 \alpha_1 \alpha_2 \mathbf{x}_1^T \mathbf{x}_2 - y_1 \alpha_1 \sum_{j=3}^n \alpha_j y_j \mathbf{x}_1^T \mathbf{x}_j \\
&\quad - y_2 \alpha_2 \sum_{j=3}^n \alpha_j y_j \mathbf{x}_2^T \mathbf{x}_j + \sum_{i=3}^n \alpha_i - \frac{1}{2} \sum_{i=3}^n \sum_{j=3}^n y_i y_j \alpha_i \alpha_j \mathbf{x}_i^T \mathbf{x}_j
\end{aligned}$$

为了描述方便，我们定义如下符号：

$$\begin{aligned}
f(\mathbf{x}_i) &= \sum_{j=1}^n \alpha_j y_j \mathbf{x}_i^T \mathbf{x}_j + b \\
v_i &= \sum_{j=3}^n \alpha_j y_j \mathbf{x}_i^T \mathbf{x}_j = f(\mathbf{x}_i) - \sum_{j=1}^2 \alpha_j y_j \mathbf{x}_i^T \mathbf{x}_j - b
\end{aligned}$$

最终目标函数变为：

$$\begin{aligned}
W(a_2) &= \alpha_1 + \alpha_2 - \frac{1}{2} \alpha_1^2 \mathbf{x}_1^T \mathbf{x}_1 - \frac{1}{2} \alpha_2^2 \mathbf{x}_2^T \mathbf{x}_2 - y_1 y_2 \alpha_1 \alpha_2 \mathbf{x}_1^T \mathbf{x}_2 - y_1 \alpha_1 v_1 - y_2 \alpha_2 v_2 \\
&\quad + \text{constant}
\end{aligned}$$

我们不关心constant的部分，因为对于 α_1 和 α_2 来说，它们都是常数项，在求导的时候，直接变为0。对于这个目标函数，如果对其求导，还有个未知推导出 α_1 和 α_2 的关系，然后用 α_2 代替 α_1 ，这样目标函数就剩一个未知数了，我们就可以求导了，推导出迭代公式。所以现在继续推导 α_1 和 α_2 的关系。注条件：

$$\sum_{i=1}^n \alpha_i y_i = 0$$

我们在求 α_1 和 α_2 的时候，可以将 $\alpha_3, \alpha_4, \dots, \alpha_n$ 和 y_3, y_4, \dots, y_n 看作常数项。因此有：

$$\alpha_1 y_1 + \alpha_2 y_2 = B$$

我们不必关心常数B的大小，现在将上述等式两边同时乘以 y_1 ，得到($y_1 y_1 = 1$)：

$$\alpha_1 = \gamma - s \alpha_2$$

其中 γ 为常数 $B y_1$ ，我们不关心这个值， $s = y_1 y_2$ 。接下来，我们将得到的 α_1 带入 $W(a_2)$ 公式得：

$$\begin{aligned}
W(a_2) &= \gamma - s \alpha_2 + \alpha_2 - \frac{1}{2} (\gamma - s \alpha_2)^2 \mathbf{x}_1^T \mathbf{x}_1 - \frac{1}{2} \alpha_2^2 \mathbf{x}_2^T \mathbf{x}_2 - s (\gamma - s \alpha_2) \alpha_2 \mathbf{x}_1^T \mathbf{x}_2 \\
&\quad - y_1 (\gamma - s \alpha_2) v_1 - y_2 \alpha_2 v_2 + \text{constant}
\end{aligned}$$

这样目标函数中就只剩下 α_2 了，我们对其求偏导（注意： $s = y_1 y_2$ ，所以 s 的平方为1， y_1 的平方和 y_2 的平方均为1）：

$$\begin{aligned}
\frac{\partial W(a_2)}{\partial \alpha_2} &= -s + 1 + \gamma s \mathbf{x}_1^T \mathbf{x}_1 - \alpha_2 \mathbf{x}_1^T \mathbf{x}_1 - \alpha_2 \mathbf{x}_2^T \mathbf{x}_2 \\
&\quad - \gamma s \mathbf{x}_1^T \mathbf{x}_2 + 2 \alpha_2 \mathbf{x}_1^T \mathbf{x}_2 + y_2 v_1 - y_2 v_2 = 0
\end{aligned}$$

继续化简，将 $s = y_1 y_2$ 带入方程。

$$\alpha_2^{\text{new}} = \frac{y_2 (y_2 - y_1 + y_1 \gamma (\mathbf{x}_1^T \mathbf{x}_1 - \mathbf{x}_1^T \mathbf{x}_2) + v_1 - v_2)}{\mathbf{x}_1^T \mathbf{x}_1 + \mathbf{x}_2^T \mathbf{x}_2 - 2 \mathbf{x}_1^T \mathbf{x}_2}$$

我们令：

$$E_i = f(x_i) - y_i$$

$$\eta = \mathbf{x}_1^T \mathbf{x}_1 + \mathbf{x}_2^T \mathbf{x}_2 - 2\mathbf{x}_1^T \mathbf{x}_2$$

E_i 为误差项， η 为学习速率。

再根据我们已知的公式：

$$\gamma = \alpha_1^{\text{old}} + s\alpha_2^{\text{old}}$$

$$v_j = \sum_{i=3}^n \alpha_i y_i \mathbf{x}_j^T \mathbf{x}_i = f(\mathbf{x}_j) - \sum_{i=1}^2 \alpha_i y_i \mathbf{x}_j^T \mathbf{x}_i - b$$

将 α_2^{new} 继续化简得：

$$\alpha_2^{\text{new}} = \alpha_2^{\text{old}} + \frac{y_2(E_1 - E_2)}{\eta}$$

这样，我们就得到了最终需要的迭代公式。这个是没有经过剪辑的解，需要考虑约束：

$$0 < \alpha_i < C$$

根据之前推导的 α 取值范围，我们得到最终的解析解为：

$$\alpha_2^{\text{new,clipped}} = \begin{cases} H & \alpha_2^{\text{new}} > H \\ \alpha_2^{\text{new}} & L \leq \alpha_2^{\text{new}} \leq H \\ L & \alpha_2^{\text{new}} < L \end{cases}$$

又因为：

$$\alpha_1^{\text{old}} = \gamma - s\alpha_2^{\text{old}}$$

$$\alpha_1^{\text{new}} = \gamma - s\alpha_2^{\text{new,clipped}}$$

消去 γ 得：

$$\alpha_1^{\text{new}} = \alpha_1^{\text{old}} + y_1 y_2 (\alpha_2^{\text{old}} - \alpha_2^{\text{new,clipped}})$$

这样，我们就知道了怎样计算 α_1 和 α_2 了，也就是如何对选择的 α 进行更新。

当我们更新了 α_1 和 α_2 之后，需要重新计算阈值 b ，因为 b 关系到了我们 $f(x)$ 的计算，也就关系到了误差 E_i 的计算。

我们要根据 α 的取值范围，去更正 b 的值，使间隔最大化。当 α_1^{new} 在0和C之间的时候，根据KKT条件可知，这个点是支持向量上的点。因此，满足

$$y_1(\omega^T \mathbf{x}_1 + b) = 1$$

公式两边同时乘以 y_1 得($y_1 y_1 = 1$)：

$$\sum_{i=1}^n \alpha_i y_i \mathbf{x}_i^T \mathbf{x}_1 + b = y_1$$

因为我们是根据 α_1 和 α_2 的值去更新 b ，所以单独提出 $i=1$ 和 $i=2$ 的时候，整理可得：

$$b_1^{\text{new}} = y_1 - \sum_{i=3}^n \alpha_i y_i \mathbf{x}_i^T \mathbf{x}_1 - \alpha_1^{\text{new}} y_1 \mathbf{x}_1^T \mathbf{x}_1 - \alpha_2^{\text{new}} y_2 \mathbf{x}_2^T \mathbf{x}_1$$

其中前两项为：

$$y_1 - \sum_{i=3}^n \alpha_i y_i \mathbf{x}_i^T \mathbf{x}_1 = -E_1 + \alpha_1^{\text{old}} y_1 \mathbf{x}_1^T \mathbf{x}_1 + \alpha_2^{\text{old}} y_2 \mathbf{x}_2^T \mathbf{x}_1 + b^{\text{old}}$$

将上述两个公式，整理得：

$$b_1^{\text{new}} = b^{\text{old}} - E_1 - y_1(\alpha_1^{\text{new}} - \alpha_1^{\text{old}}) \mathbf{x}_1^T \mathbf{x}_1$$

$$- y_2(\alpha_2^{\text{new}} - \alpha_2^{\text{old}}) \mathbf{x}_2^T \mathbf{x}_1$$

同理可得 b_2 new为：

$$b_2^{new} = b^{old} - E_2 - y_1(\alpha_1^{new} - \alpha_1^{old})\mathbf{x}_1^T \mathbf{x}_2 - y_2(\alpha_2^{new} - \alpha_2^{old})\mathbf{x}_2^T \mathbf{x}_2$$

当 b_1 和 b_2 都有效的时候，它们是相等的，即：

$$b^{new} = b_1^{new} = b_2^{new}$$

当两个乘子都在边界上，则 b 阈值和KKT条件一致。当不满足的时候，SMO算法选择他们的中点作为新的阈值：

$$b = \begin{cases} b_1, & 0 < \alpha_1^{new} < C \\ b_2 & 0 < \alpha_2^{new} < C \\ (b_1 + b_2)/2 & \text{otherwise} \end{cases}$$

最后，更新所有的 α 和 b ，这样模型就出来了，从而即可求出我们的分类函数。

现在，让我们梳理下SMO算法的步骤：

■ 步骤1：计算误差：

$$E_i = f(\mathbf{x}_i) - y_i = \sum_{j=1}^n \alpha_j y_j \mathbf{x}_i^T \mathbf{x}_j + b - y_i$$

■ 步骤2：计算上下界 L 和 H ：

$$\begin{cases} L = \max(0, \alpha_j^{old} - \alpha_i^{old}), H = \min(C, C + \alpha_j^{old} - \alpha_i^{old}) & \text{if } y_i \neq y_j \\ L = \max(0, \alpha_j^{old} + \alpha_i^{old} - C), H = \min(C, \alpha_j^{old} + \alpha_i^{old}) & \text{if } y_i = y_j \end{cases}$$

■ 步骤3：计算 η ：

$$\eta = \mathbf{x}_i^T \mathbf{x}_i + \mathbf{x}_j^T \mathbf{x}_j - 2\mathbf{x}_i^T \mathbf{x}_j$$

■ 步骤4：更新 α_j ：

$$\alpha_j^{new} = \alpha_j^{old} + \frac{y_j(E_i - E_j)}{\eta}$$

■ 步骤5：根据取值范围修剪 α_j ：

$$\alpha_j^{new,clipped} = \begin{cases} H & \text{if } \alpha_j^{new} \geq H \\ \alpha_j^{new} & \text{if } L \leq \alpha_j^{new} \leq H \\ L & \text{if } \alpha_j^{new} \leq L \end{cases}$$

■ 步骤6：更新 α_i ：

$$\alpha_i^{new} = \alpha_i^{old} + y_i y_j (\alpha_j^{old} - \alpha_j^{new,clipped})$$

■ 步骤7：更新 b_1 和 b_2 ：

$$\begin{aligned} b_1^{new} &= b^{old} - E_i - y_i(\alpha_i^{new} - \alpha_i^{old})\mathbf{x}_i^T \mathbf{x}_i - y_j(\alpha_j^{new} - \alpha_j^{old})\mathbf{x}_j^T \mathbf{x}_i \\ b_2^{new} &= b^{old} - E_j - y_i(\alpha_i^{new} - \alpha_i^{old})\mathbf{x}_i^T \mathbf{x}_j - y_j(\alpha_j^{new} - \alpha_j^{old})\mathbf{x}_j^T \mathbf{x}_j \end{aligned}$$

■ 步骤8：根据 b_1 和 b_2 更新 b ：

$$b = \begin{cases} b_1 & 0 < \alpha_1^{new} < C \\ b_2 & 0 < \alpha_2^{new} < C \\ \frac{b_1 + b_2}{2} & \text{otherwise} \end{cases}$$

四、编程求解线性SVM

已经梳理完了SMO算法实现步骤，接下来按照这个思路编写代码，进行实战练习。

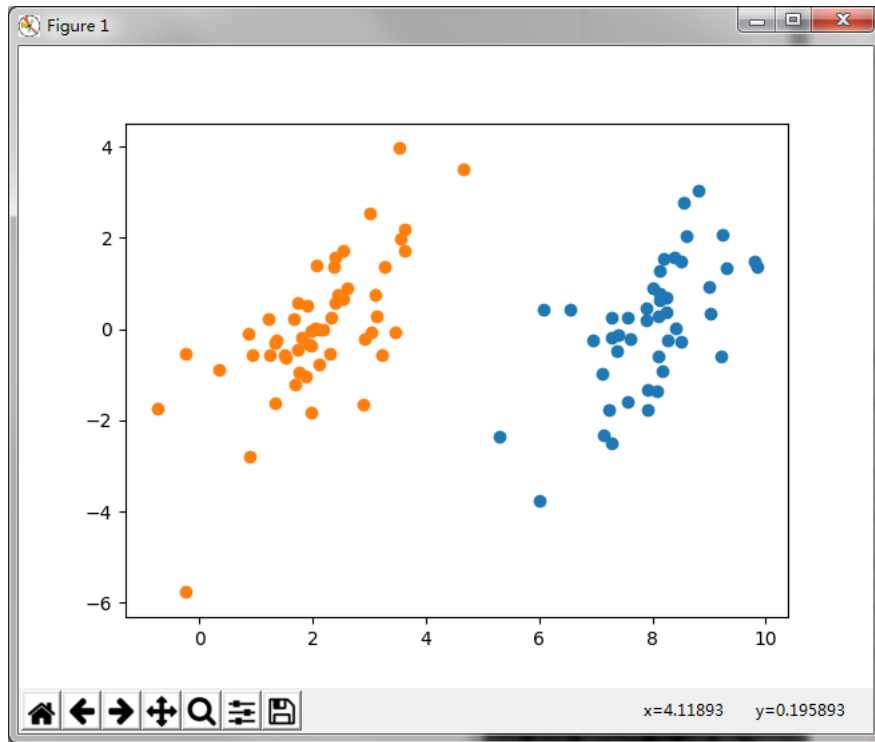
1、可视化数据集

我们先使用简单的数据集进行测试，数据集下载地址：[数据集下载](#)

编写程序可视化数据集，看下它是长什么样的：

```
1  # -*- coding:UTF-8 -*-
2  import matplotlib.pyplot as plt
3  import numpy as np
4
5  """
6  函数说明:读取数据
7
8  Parameters:
9      fileName - 文件名
10 Returns:
11     dataMat - 数据矩阵
12     labelMat - 数据标签
13 Author:
14     Jack Cui
15 Blog:
16     http://blog.csdn.net/c406495762
17 Zhihu:
18     https://www.zhihu.com/people/Jack--Cui/
19 Modify:
20     2017-09-21
21 """
22 def loadDataSet(fileName):
23     dataMat = []; labelMat = []
24     fr = open(fileName)
25     for line in fr.readlines():
26         lineArr = line.strip().split('\t')
27         dataMat.append([float(lineArr[0]), float(lineArr[1])])
28         labelMat.append(float(lineArr[2]))
29     return dataMat, labelMat
30
31 """
32 函数说明:数据可视化
33
34 Parameters:
35     dataMat - 数据矩阵
36     labelMat - 数据标签
37 Returns:
38     无
39 Author:
40     Jack Cui
41 Blog:
42     http://blog.csdn.net/c406495762
43 Zhihu:
44     https://www.zhihu.com/people/Jack--Cui/
45 Modify:
46     2017-09-21
47 """
48 def showDataSet(dataMat, labelMat):
49     data_plus = []
50     data_minus = []
51     for i in range(len(dataMat)):
52         if labelMat[i] > 0:
53             data_plus.append(dataMat[i])
54         else:
55             data_minus.append(dataMat[i])
56     data_plus_np = np.array(data_plus)
57     data_minus_np = np.array(data_minus)
58     plt.scatter(np.transpose(data_plus_np)[0], np.transpose(data_plus_np)[1])
59     plt.scatter(np.transpose(data_minus_np)[0], np.transpose(data_minus_np)[1])
60     plt.show()
61
62 if __name__ == '__main__':
63     dataMat, labelMat = loadDataSet('testSet.txt')
64     showDataSet(dataMat, labelMat)
```

运行程序，查看结果：



这就是我们使用的二维数据集，显然线性可分。现在我们使用简化版的SMO算法进行求解。

2、简化版SMO算法

按照上述已经推导的步骤编写代码：

```

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  Author:
17      Jack Cui
18  Blog:
19      http://blog.csdn.net/c406495762
20  Zhihu:
21      https://www.zhihu.com/people/Jack--Cui/
22  Modify:
23      2017-09-21
24  """
25  def loadDataSet(fileName):
26      dataMat = []; labelMat = []
27      fr = open(fileName)
28      for line in fr.readlines():
29          lineArr = line.strip().split('\t')
30          dataMat.append([float(lineArr[0]), float(lineArr[1])])
31          labelMat.append(float(lineArr[2]))
32      return dataMat, labelMat
33
34  """
35  函数说明:随机选择alpha
36
37  Parameters:
38      i - alpha
39      m - alpha参数个数
40  Returns:
41      j -
42  Author:
43      Jack Cui
44  Blog:
45      http://blog.csdn.net/c406495762
46  Zhihu:
47      https://www.zhihu.com/people/Jack--Cui/
48  Modify:
49      2017-09-21
50  """
51  def selectJrand(i, m):
52      j = i
53      while (j == i):
54          j = int(random.uniform(0, m))
55      return j
56
57  """
58  """

```

```

59 函数说明:修剪alpha
60
61 Parameters:
62     aj - alpha值
63     H - alpha上限
64     L - alpha下限
65 Returns:
66     aj - alpah值
67 Author:
68     Jack Cui
69 Blog:
70     http://blog.csdn.net/c406495762
71 Zhihu:
72     https://www.zhihu.com/people/Jack--Cui/
73 Modify:
74     2017-09-21
75 """
76 def clipAlpha(aj,H,L):
77     if aj > H:
78         aj = H
79     if L > aj:
80         aj = L
81     return aj
82
83 """
84 函数说明:简化版SMO算法
85
86 Parameters:
87     dataMatIn - 数据矩阵
88     classLabels - 数据标签
89     C - 松弛变量
90     toler - 容错率
91     maxIter - 最大迭代次数
92 Returns:
93     无
94 Author:
95     Jack Cui
96 Blog:
97     http://blog.csdn.net/c406495762
98 Zhihu:
99     https://www.zhihu.com/people/Jack--Cui/
100 Modify:
101     2017-09-23
102 """
103 def smoSimple(dataMatIn, classLabels, C, toler, maxIter):
104     #转换为numpy的mat存储
105     dataMatrix = np.mat(dataMatIn); labelMat = np.mat(classLabels).transpose()
106     #初始化b参数, 统计dataMatrix的维度
107     b = 0; m,n = np.shape(dataMatrix)
108     #初始化alpha参数, 设为0
109     alphas = np.mat(np.zeros((m,1)))
110     #初始化迭代次数
111     iter_num = 0
112     #最多迭代maxIter次
113     while (iter_num < maxIter):
114         alphaPairsChanged = 0
115         for i in range(m):
116             #步骤1: 计算误差Ei
117             fXi = float(np.multiply(alphas,labelMat).T*(dataMatrix*dataMatrix[i,:].T)) + b
118             Ei = fXi - float(labelMat[i])
119             #优化alpha, 更设定一定的容错率。
120             if ((labelMat[i]*Ei < -toler) and (alphas[i] < C)) or ((labelMat[i]*Ei > toler) and (alphas[i] > 0)):
121                 #随机选择另一个与alpha_i成对优化的alpha_j
122                 j = selectJrand(i,m)
123                 #步骤1: 计算误差Ej
124                 fXj = float(np.multiply(alphas,labelMat).T*(dataMatrix*dataMatrix[j,:].T)) + b
125                 Ej = fXj - float(labelMat[j])
126                 #保存更新前的alpha值, 使用深拷贝
127                 alphaJold = alphas[j].copy(); alphaJold = alphas[j].copy();
128                 #步骤2: 计算上下界L和H
129                 if (labelMat[i] != labelMat[j]):
130                     L = max(0, alphas[j] - alphas[i])
131                     H = min(C, C + alphas[j] - alphas[i])
132                 else:
133                     L = max(0, alphas[j] + alphas[i] - C)
134                     H = min(C, alphas[j] + alphas[i])
135                 if L==H: print("L==H"); continue
136                 #步骤3: 计算eta
137                 eta = 2.0 * dataMatrix[i,:]*dataMatrix[j,:].T - dataMatrix[i,:]*dataMatrix[i,:].T - dataMatrix[j,:]*dataMatrix[j,:].T
138                 if eta >= 0: print("eta>=0"); continue
139                 #步骤4: 更新alpha_j
140                 alphas[j] -= labelMat[j]*(Ei - Ej)/eta
141                 #步骤5: 修剪alpha_j
142                 alphas[j] = clipAlpha(alphas[j],H,L)
143                 if (abs(alphas[j] - alphaJold) < 0.00001): print("alpha_j变化太小"); continue
144                 #步骤6: 更新alpha_i
145                 alphas[i] += labelMat[j]*labelMat[i]*(alphaJold - alphas[j])
146                 #步骤7: 更新b_1和b_2
147                 b1 = b - Ei - labelMat[i]*(alphas[i]-alphaJold)*dataMatrix[i,:]*dataMatrix[i,:].T - labelMat[j]*(alphas[j]-alphaJold)*dataMatrix[
148                 b2 = b - Ej - labelMat[i]*(alphas[i]-alphaJold)*dataMatrix[i,:]*dataMatrix[j,:].T - labelMat[j]*(alphas[j]-alphaJold)*dataMatrix[
149                 #步骤8: 根据b_1和b_2更新b
150                 if (0 < alphas[i]) and (C > alphas[i]): b = b1
151                 elif (0 < alphas[j]) and (C > alphas[j]): b = b2
152                 else: b = (b1 + b2)/2.0
153                 #统计优化次数
154                 alphaPairsChanged += 1
155                 #打印统计信息
156                 print("第%d次迭代 样本:%d, alpha优化次数:%d" % (iter_num,i,alphaPairsChanged))
157             #更新迭代次数
158             if (alphaPairsChanged == 0): iter_num += 1
159             else: iter_num = 0
160             print("迭代次数: %d" % iter_num)
161     return b,alphas
162

```

```

163 """
164 函数说明:分类结果可视化
165
166 Parameters:
167     dataMat - 数据矩阵
168     w - 直线法向量
169     b - 直线解决
170 Returns:
171     无
172 Author:
173     Jack Cui
174 Blog:
175     http://blog.csdn.net/c406495762
176 Zhihu:
177     https://www.zhihu.com/people/Jack--Cui/
178 Modify:
179     2017-09-23
180 """
181 def showClassifier(dataMat, w, b):
182     #绘制样本点
183     data_plus = [] #正样本
184     data_minus = [] #负样本
185     for i in range(len(dataMat)):
186         if labelMat[i] > 0:
187             data_plus.append(dataMat[i])
188         else:
189             data_minus.append(dataMat[i])
190     data_plus_np = np.array(data_plus) #转换为numpy矩阵
191     data_minus_np = np.array(data_minus) #转换为numpy矩阵
192     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
193     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
194     #绘制直线
195     x1 = max(dataMat)[0]
196     x2 = min(dataMat)[0]
197     a1, a2 = w
198     b = float(b)
199     a1 = float(a1[0])
200     a2 = float(a2[0])
201     y1, y2 = (-b- a1*x1)/a2, (-b - a1*x2)/a2
202     plt.plot([x1, x2], [y1, y2])
203     #找出支持向量点
204     for i, alpha in enumerate(alphas):
205         if alpha > 0:
206             x, y = dataMat[i]
207             plt.scatter([x], [y], s=150, c='none', alpha=0.7, linewidth=1.5, edgecolor='red')
208     plt.show()
209
210 """
211 函数说明:计算w
212
213 Parameters:
214     dataMat - 数据矩阵
215     labelMat - 数据标签
216     alphas - alphas值
217 Returns:
218     无
219 Author:
220     Jack Cui
221 Blog:
222     http://blog.csdn.net/c406495762
223 Zhihu:
224     https://www.zhihu.com/people/Jack--Cui/
225 Modify:
226     2017-09-23
227 """
228
229 def get_w(dataMat, labelMat, alphas):
230     alphas, dataMat, labelMat = np.array(alphas), np.array(dataMat), np.array(labelMat)
231     w = np.dot((np.tile(labelMat.reshape(1, -1).T, (1, 2)) * dataMat).T, alphas)
232     return w.tolist()
233
234
235 if __name__ == '__main__':
236     dataMat, labelMat = loadDataSet('testSet.txt')
237     b, alphas = smoSimple(dataMat, labelMat, 0.6, 0.001, 40)
238     w = get_w(dataMat, labelMat, alphas)
239     showClassifier(dataMat, w, b)

```

程序运行结果：

其中，中间的蓝线为求出来的分类器，用红圈圈出的点为支持向量点。

五、总结

- 本文主要进行了线性SVM的推导，并通过编程实现一个简化版SMO算法；
- 本文的简化版SMO算法在选取 α 的时候，没有选择启发式的选择方法，并且两个乘子的计算没有进行优化，所以算法比较耗时，下一篇文章会讲解优化方法；
- 本文讨论的是线性SVM，没有使用核函数，下一篇文章将会讲解如何应用核函数，将SVM应用于非线性数据集；
- 如有问题，请留言。如有错误，还望指正，谢谢！

PS：如果觉得本篇本章对您有所帮助，欢迎关注、评论、赞！

参考资料：

- [1] 五岁小孩也能看懂的SVM：<https://www.zhihu.com/question/21094489/answer/86273196>
- [2] 五岁小孩也能看懂的SVM：
https://www.reddit.com/r/MachineLearning/comments/15zrpp/please_explain_support_vector_machines_svm_like_i/
- [3] pluskid大牛博客：http://blog.pluskid.org/?page_id=683
- [4] 陈东岳老师文章：<https://zhuanlan.zhihu.com/p/24638007>
- [5] 深入理解拉格朗日乘子法和KKT条件：<http://blog.csdn.net/xianlingmao/article/details/7919597>
- [6] 充分条件和必要条件：<https://www.zhihu.com/question/30469121>
- [7] 凸函数：<https://zh.wikipedia.org/wiki/%E5%87%B8%E5%87%BD%E6%95%B0>
- [8] 《机器学习实战》第6章内容。
- [9] SVM之SMO算法：<http://www.cnblogs.com/zangrunqiang/p/5515872.html>



JackCui

关注人工智能及互联网的个人博客

[查看熊掌号](#)