

支持向量机(Support Vector Machines, SVM)

线性可分支持向量机

线性可分支持向量机

给定训练数据：

$$T = \{(x_1, y_1), (x_2, y_2), \dots, (x_N, y_N)\} \quad (1)$$

其中 $x_i \in R^n, y_i \in \{1, -1\}, i = 1, 2, \dots, N$ 可以被一超平面 $w \cdot x + b = 0$ 根据类别 y_i 分开。

给定线性可分训练数据集，通过间隔最大化或等价地求解相应的凸二次规划问题学习得到的分离超平面为

$$w^* \cdot x + b^* = 0$$

以及相应的分类决策函数

$$f(x) = \text{sign}(w^* \cdot x + b^*)$$

称为线性可分支持向量机

函数间隔和几何间隔

对于给定的训练集 T 以及超平面 (w, b) ，任意的 (x_i, y_i) 到超平面的距离为：

$$l = \frac{|w \cdot x_i + b|}{\|w\|} \quad (2)$$

其中分子 $|w \cdot x_i + b|$ 可以相对的表示 (x_i, y_i) 到超平面的距离，可以将其认为分类的确信程度(离超平面越远确信程度越高)，而 y_i 和 $w \cdot x_i + b$ 的乘积可以表示分类是否正确，所以可以用 $y_i(w \cdot x_i + b)$ 来表示分类的正确性以及确信度。

函数间隔

对于给定的训练集 T 以及超平面 (w, b) ，定义超平面关于样本点 (x_i, y_i) 的函数间隔为：

$$\hat{\gamma}_i = y_i(w \cdot x_i + b)$$

而超平面关于训练集 T 的函数间隔定义为：

$$\hat{\gamma} = \min_{i=1,2,\dots,N} \hat{\gamma}_i$$

几何间隔

注意到，当 w, b 成倍数增长时，函数间隔 γ 也会成倍数增长，即：

$$\begin{aligned} w &= \lambda w \\ b &= \lambda b \\ \gamma &\rightarrow \lambda \gamma \end{aligned} \quad (3)$$

但是超平面 (w, b) 是没有发生变化的：

$$\lambda w \cdot x + \lambda b = 0 \quad (4)$$

为了唯一的确定超平面，我们可以对 w 做一些约束，使得 $\|w\| = 1$ ，从而可以得到唯一的 w, b 。

定义超平面对样本点 (x_i, y_i) 的几何间隔为：

$$\gamma_i = y_i \left(\frac{w}{\|w\|} \cdot x + \frac{b}{\|w\|} \right)$$

定义超平面对训练集的几何间隔为：

$$\gamma = \min_{i=1,2,\dots,N} \gamma_i$$

根据函数间隔和几何间隔的定义可以得出两者的关系：

$$\begin{aligned} \gamma_i &= \frac{\hat{\gamma}_i}{\|w\|} \\ \gamma &= \frac{\hat{\gamma}}{\|w\|} \end{aligned} \quad (5)$$

间隔最大化

对线性可分的训练数据集而言，线性可分分离超平面有无穷多个(等价于感知机)，但是几何间隔最大的分离超平面却是唯一的。这里的间隔最大化又称为硬间隔最大化。

间隔最大化的直观解释是：对训练数据集找到几何间隔最大的超平面意味着以充分大的确信度对训练数据进行分类。也就是说，不仅将正负实例点分开，而且对最难分的实例点(离超平面最近的点)也有足够大的确信度将他们分开。这样的超平面应该对未知的新实例有很好的分类预测能力。

最大间隔分离超平面

几何间隔最大的分离超平面，并且满足一些约束条件，等价于下面的约束最优化问题：

$$\begin{aligned} \max_{w,b} \quad & \gamma \\ \text{s.t.} \quad & y_i \left(\frac{w \cdot x}{\|w\|} + \frac{b}{\|w\|} \right) \geq \gamma \end{aligned} \quad (6)$$

考虑函数间隔和几何间隔的关系，这个问题也等价于：

$$\begin{aligned} \max_{w,b} \quad & \frac{\hat{\gamma}}{\|w\|} \\ \text{s.t.} \quad & y_i (w \cdot x + b) \geq \hat{\gamma} \end{aligned} \quad (7)$$

函数间隔的取值并不影响最优化问题的解(不影响模型)，也就是说，当 w, b 成倍数增长时，函数间隔也成相同倍数增长，但是模型(超平面)是不变的。反之，如果函数间隔不确定， w, b 会有无穷多个解(都成一定的倍数关系)，为了确定出唯一解，这里我们设定 $\hat{\gamma} = 1$ ，那么问题就转化为：

$$\begin{aligned} \max_{w,b} \quad & \frac{1}{\|w\|} \\ \text{s.t.} \quad & y_i (w \cdot x + b) \geq 1 \end{aligned} \quad (8)$$

为了方便后续计算，上述问题等价于：

$$\begin{aligned} \min_{w,b} \quad & \|w\| \\ \text{s.t.} \quad & y_i (w \cdot x + b) \geq 1 \end{aligned} \quad (9)$$

为了使构造函数是凸函数，将上述问题继续等价：

$$\begin{aligned} \min_{w,b} \quad & \frac{1}{2} \|w\|^2 \\ \text{s.t.} \quad & 1 - y_i (w \cdot x + b) \leq 0 \end{aligned} \quad (10)$$

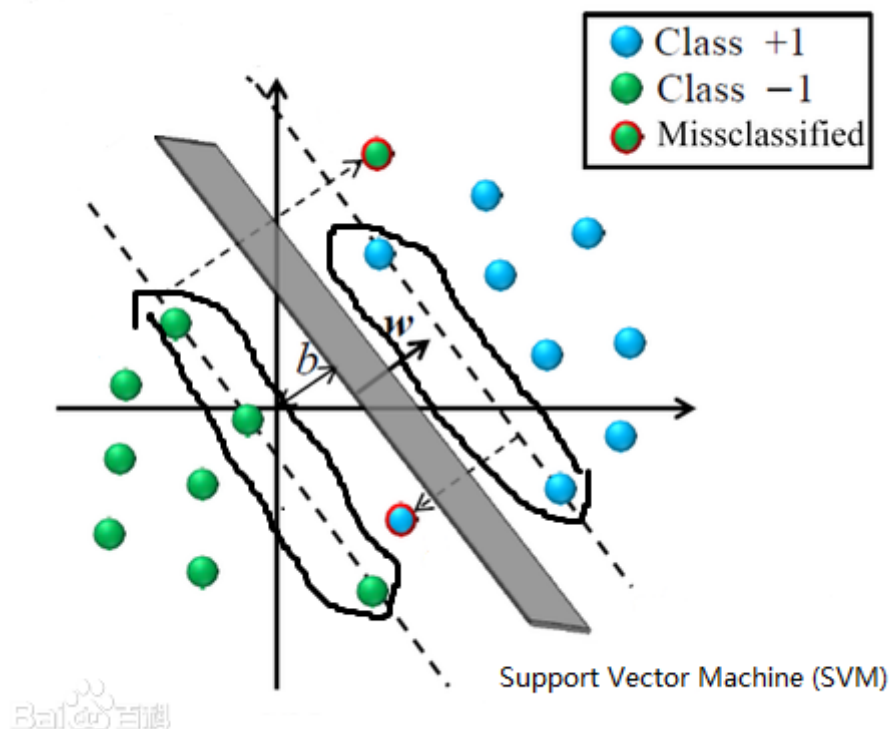
得到上述的等价问题就可以应用拉格朗日数乘法了。

支持向量和间隔边界

在线性可分的情况下，训练数据集的样本点中与分离超平面距离最近的样本点的实例称为支持向量 (support vector)。支持向量是使约束条件 $1 - y_i(w \cdot x + b) \leq 0$ 成立的点，也即 $y_i(w \cdot x + b) = 1$ ：

对 $y_i = +1$ 的正例点，支持向量在超平面 $H_1 : w \cdot x + b = 1$ 上

对 $y_i = -1$ 的负例点，支持向量在超平面 $H_2 : w \cdot x + b = -1$ 上



H_1 和 H_2 即为间隔边界，而 H_1 和 H_2 之间的距离称为间隔 (margin)。间隔依赖于分离超平面的法向量 w ，等于 $\frac{2}{\|w\|}$ 。

在决定分离超平面时，实际上只有支持向量起作用，所以这种分类模型称之为“支持向量机 (SVM)”。

学习的对偶算法

为了求解出超平面 (w, b) ，我们需要满足约束条件以及求解目标，原始问题为：

$$\begin{aligned} \min_{w,b} \quad & \frac{1}{2} \|w\|^2 \\ \text{s.t.} \quad & 1 - y_i(w \cdot x + b) \leq 0 \end{aligned} \quad (11)$$

构建拉格朗日函数 $L(w, b, a)$ ：

$$\begin{aligned} L(w, b, a) &= \frac{1}{2} \|w\|^2 + \sum_{i=1}^N a_i (1 - y_i(w \cdot x_i + b)) \\ &= \frac{1}{2} \|w\|^2 - \sum_{i=1}^N a_i y_i (w \cdot x_i + b) + \sum_{i=1}^N a_i \end{aligned}$$

原始问题为：

$$\min_{w,b} \max_a L(w, b, a) \quad (12)$$

对偶问题为：

$$\max_a \min_{w,b} L(w, b, a) \quad (13)$$

接下来对对偶问题进行求解，从里到外进行求解

1. 先求 $\min_{w,b} L(w, b, a)$

对 w, b 求偏导并令其等于0。

$$\nabla_w L(w, b, a) = w - \sum_{i=1}^N a_i y_i x_i \quad (14)$$

$$\nabla_b L(w, b, a) = - \sum_{i=1}^N a_i y_i$$

(其中 $\|w\|^2 = w^T w = w^T I w$, (I 为单位矩阵), 根据矩阵求导法则 $\frac{\partial (X^T A X)}{\partial X} = (A + A^T) X$ 可知, 其求导结果为 $2Iw = 2w$, 而根据 $\frac{\partial A X}{\partial X} = A$ 可求出其余项的偏导)

$$\begin{aligned} w &= \sum_{i=1}^N a_i y_i x_i \\ \sum_{i=1}^N a_i y_i &= 0 \end{aligned} \quad (15)$$

将其带回到 $L(w, b, a)$ 得到:

$$\begin{aligned} \sum_{i=1}^N a_i y_i (w \cdot x_i + b) + \sum_{i=1}^N a_i \\ \min_{w,b} L(w, b, a) &= \frac{1}{2} \sum_{i,j=1}^N a_i a_j y_i y_j (x_i \cdot x_j) - \sum_{i=1}^N a_i y_i \left(\sum_{j=1}^N a_j y_j (x_j \cdot x_i) + b \right) + \sum_{i=1}^N a_i \\ &= \frac{1}{2} \sum_{i,j=1}^N a_i a_j y_i y_j (x_i \cdot x_j) - \sum_{i,j=1}^N a_i a_j y_i y_j (x_i \cdot x_j) - b \sum_{i=1}^N a_i y_i + \sum_{i=1}^N a_i \\ &= -\frac{1}{2} \sum_{i,j=1}^N a_i a_j y_i y_j (x_i \cdot x_j) + \sum_{i=1}^N a_i \end{aligned}$$

2. 求 $\min_{w,b} L(w, b, a)$ 的极大, 即是对偶问题:

$$\begin{aligned} \max_a \quad & -\frac{1}{2} \sum_{i,j=1}^N a_i a_j y_i y_j (x_i \cdot x_j) + \sum_{i=1}^N a_i \quad (16) \\ \text{s.t.} \quad & \sum_{i=1}^N a_i y_i = 0 \\ & a_i \geq 0, \quad i = 1, 2, \dots, N \end{aligned}$$

线性支持向量机和软间隔最大化

线性支持向量机

线性可分问题的支持向量机学习方法, 对线性不可分训练数据是不适用的, 因为这时候原不等式约束 $y_i(w \cdot x + b) \geq 1$ 并不能都成立。所以引入了软间隔最大化。

给定训练数据:

$$T = \{(x_1, y_1), (x_2, y_2), \dots, (x_N, y_N)\} \quad (17)$$

其中 $x_i \in R^n, y_i \in \{1, -1\}, i = 1, 2, \dots, N$ 。假定这些训练数据集**不是线性可分**的。一般情况下, 训练数据中可能会混有一些**特异点(outlier)**, 而将这些特异点除去后, 剩下的大部分样本点组成的集合是线性可分的。

为了解决特异点造成的线性不可分问题, 可以在原来约束的基础上引入**松弛变量** $\xi_i \geq 0$, 使得函数间隔加上 ξ_i 能大于等于1, 这样的话, 约束条件就变成了:

$$\begin{aligned}y_i(w \cdot x_i + b) + \xi_i &\geq 1 \\y_i(w \cdot x_i + b) &\geq 1 - \xi_i\end{aligned}$$

当出现特异点的时候 $\xi_i > 0$ ，特别地当出现误分类点时 $\xi_i > 1$ 。但是如果 ξ_i 很大很大时，随便一个超平面都满足上面的条件，所以我们需要尽量减小 ξ_i 的值来确定正确的超平面。

所以将目标函数设置为：

$$\min_{w,b,\xi} \frac{1}{2} \|w\|^2 + C \sum_{i=1}^N \xi_i \quad (18)$$

其中 $C > 0$ 为惩罚参数，是一个调和 w 、 ξ_i 的权重系数，当 C 比较大时，为了保证整体比较小，那么 ξ_i 尽可能小的取值也就是尽可能取0，也就是说此时要求模型的**特异点尽可能少**。

原始问题

$$\begin{aligned}\min_{w,b,\xi} \quad & \frac{1}{2} \|w\|^2 + C \sum_{i=1}^N \xi_i \\s.t. \quad & y_i(w \cdot x_i + b) \geq 1 - \xi_i \\& \xi_i \geq 0 \quad i = 1, 2, \dots, N\end{aligned}$$

线性支持向量机

对于给定的线性不可分的训练数据集，通过求解上述原始的二次凸优化问题，也即软间隔最大化问题，得到的分离超平面为：

$$w^* \cdot x + b^* = 0 \quad (19)$$

以及相应的分类决策函数：

$$f(x) = \text{sign}(w^* \cdot x + b^*) \quad (20)$$

称为线性支持向量机

学习的对偶算法

构造拉格朗日函数 L ：

$$L(w, b, a, \mu, \xi) = \frac{1}{2} \|w\|^2 + C \sum_{i=1}^N \xi_i + \sum_{i=1}^N a_i (1 - \xi_i - y_i(w \cdot x_i + b)) - \sum_{i=1}^N \mu_i \xi_i \quad (21)$$

对偶问题

$$\max_{a,\mu} \min_{w,b,\xi} L(w, b, a, \mu, \xi) \quad (22)$$

1. 先求 $\min_{w,b,\xi} L(w, b, a, \mu, \xi)$ ：

$$\nabla_w L = w - \sum_{i=1}^N a_i y_i x_i = 0$$

$$\nabla_b L = - \sum_{i=1}^N a_i y_i = 0$$

$$\nabla_{\xi_i} L = C - a_i - \mu_i = 0$$

得到：

$$w = \sum_{i=1}^N a_i y_i x_i \quad (23)$$

$$\sum_{i=1}^N a_i y_i = 0$$

$$C - a_i - \mu_i = 0$$

带回到 $\min_{w,b,\xi} L(w, b, a, \mu, \xi)$ 中的到:

$$\min_{w,b,\xi} L(w, b, a, \mu, \xi) = -\frac{1}{2} \sum_{i,j=1}^N a_i a_j y_i y_j (x_i \cdot x_j) + \sum_{i=1}^N a_i \quad (24)$$

2. 求 $\min_{w,b,\xi} L(w, b, a, \mu, \xi)$ 的极大, 即求 $\max_{a,\mu} \min_{w,b,\xi} L(w, b, a, \mu, \xi)$:

$$\max_{a,\mu} -\frac{1}{2} \sum_{i,j=1}^N a_i a_j y_i y_j (x_i \cdot x_j) + \sum_{i=1}^N a_i$$

$$s. t. \quad \sum_{i=1}^N a_i y_i = 0$$

$$C - a_i - \mu_i = 0$$

$$a_i \geq 0$$

$$\mu_i \geq 0$$

因为 $\min_{w,b,\xi} L(w, b, a, \mu, \xi)$ 中 μ 以及被消掉, 整理上述式子可以得到:

$$\max_a -\frac{1}{2} \sum_{i,j=1}^N a_i a_j y_i y_j (x_i \cdot x_j) + \sum_{i=1}^N a_i$$

$$s. t. \quad \sum_{i=1}^N a_i y_i = 0$$

$$0 \leq a_i \leq C$$

支持向量

讨论一下上述对偶问题需要满足的 KKT 条件:

$$\text{原始拉格朗日函数: } L(w, b, a, \mu, \xi) = \frac{1}{2} \|w\|^2 + C \sum_{i=1}^N \xi_i + \sum_{i=1}^N a_i (1 - \xi_i - y_i (w \cdot x_i + b)) - \sum_{i=1}^N \mu_i \xi_i$$

$$\nabla_{w^*} L = w^* - \sum_{i=1}^N a_i^* y_i x_i = 0$$

$$\nabla_{b^*} L = \sum_{i=1}^N a_i^* y_i = 0$$

$$\nabla_{\xi_i^*} L = C - a_i^* - \mu_i^* = 0$$

$$a_i^* (1 - \xi_i^* - y_i (w_i^* \cdot x + b^*)) = 0$$

$$\mu_i^* \xi_i^* = 0$$

$$\xi_i^* \geq 0$$

$$1 - \xi_i^* - y_i (w_i^* \cdot x + b^*) \leq 0$$

$$a_i^* \geq 0$$

$$\mu_i^* \geq 0$$

对偶问题以及原始问题通解需要满足上述 KKT 条件, 所以:

- $a_i^* = 0$ 时, $u_i^* = C$, 必有 $\xi_i = 0$, 所以有 $y_i (w_i^* \cdot x + b^*) \geq 1$:

此时样本点在间隔边界上或者分类正确

- $a_i^* = C$ 时, $u_i^* = 0$, 必有 $\xi_i^* \geq 0$, 且有 $y_i (w_i^* \cdot x + b^*) = 1 - \xi_i^*$:

此时根据 ξ_i^* 的取值样本点的位置可分为分类错误、超平面上、超平面和间隔边界之间、间隔边界上

- 当 $0 < a_i^* < C$ 时, 必有 $\xi_i^* = 0$, 所以有 $y_i (w_i^* \cdot x + b^*) = 1$:

此时直线在间隔边界上, 所以此时的样本点 (x_i, y_i) 为支持向量

非线性支持向量机和核函数

写出核函数篇幅较长，而这里需要的核函数知识较少，读者可参考这位大佬的[文章](#)

序列最小最优算法(SMO)

数据集较大时，训练SVM需要求解一个大规模的二次规划(quadratic programming, QP)问题。SMO算法的大致思路是：将这个大的QP问题分解成一系列尽可能小的QP子问题，然后依次求解这些小的QP子问题(Osuna等人已经证明当满足一定条件时可以将大QP问题分解为小QP问题，并且可以收敛)

SMO算法的停机条件

前面谈到的对偶问题是：

$$\begin{aligned} \min_a \quad & \frac{1}{2} \sum_{i,j=1}^N a_i a_j y_i y_j K(\mathbf{x}_i, \mathbf{x}_j) - \sum_{i=1}^N a_i \quad (25) \\ \text{s.t.} \quad & 0 \leq a_i \leq C \\ & \sum_{i=1}^N a_i y_i = 0 \end{aligned}$$

$K(\mathbf{x}_i, \mathbf{x}_j)$ 为正定核函数，下面简记为 K_{ij} 。

KKT条件重温：

$$\begin{aligned} \nabla_{w^*} L &= w^* - \sum_{i=1}^N a_i^* y_i x_i = 0 \\ \nabla_{b^*} L &= \sum_{i=1}^N a_i^* y_i = 0 \\ \nabla_{\xi_i^*} L &= C - a_i^* - \mu_i^* = 0 \\ a_i^* (1 - \xi_i^* - y_i (w_i^* \cdot x + b^*)) &= 0 \\ \mu_i^* \xi_i^* &= 0 \\ \xi_i^* &\geq 0 \\ 1 - \xi_i^* - y_i (w_i^* \cdot x + b^*) &\leq 0 \\ a_i^* &\geq 0 \\ \mu_i^* &\geq 0 \end{aligned}$$

以及：

$$\begin{aligned} a_i = 0 &\Rightarrow y_i \left(\sum_{j=1}^N a_j y_j K_{ij} + b \right) \geq 1 \\ 0 < a_i < C &\Rightarrow y_i \left(\sum_{j=1}^N a_j y_j K_{ij} + b \right) = 1 \\ a_i = C &\Rightarrow y_i \left(\sum_{j=1}^N a_j y_j K_{ij} + b \right) \leq 1 \end{aligned}$$

为方便表示，令 $g(x_i) = \sum_{j=1}^N a_j y_j K_{ij} + b$,

根据Osuna以及Platt的研究，停机条件为：

$$\begin{aligned}
0 &\leq a_i \leq C \\
\sum_{i=1}^N a_i y_i &= 0 \\
y_i g(x_i) &= \begin{cases} \geq 1, \{x_i | a_i = 0\} \\ = 1, \{x_i | 0 < a_i < C\} \\ \leq 1, \{x_i | a_i = C\} \end{cases}
\end{aligned}$$

SMO算法的思路

如果训练集中的所有样本都满足上述停机条件，那么就得到了问题的最优解；否则将大QP问题拆分成一系列小的QP问题，每次对小的QP问题求解，然后一步步向最优解优化。

最小的QP子问题包含两个变量，但是这两个变量是互相受限的，因为有等式 $\sum_{i=1}^N a_i y_i = 0$ 的存在。如果每个最小QP问题只优化一个变量则无法满足上述等式约束，所以可以两个变量一起进行优化，从而可以保证等式约束是成立的。

所以现在我们的重点在于：**两个变量二次规划的解析方法** and **选择变量的启发式方法**。

两个变量二次规划的解析方法

不失一般性，假设选择的两个变量为 a_1, a_2 ，其他变量视为常量。

原目标函数为：

$$\min_a \frac{1}{2} \sum_{i,j=1}^N a_i a_j y_i y_j K(\mathbf{x}_i, \mathbf{x}_j) - \sum_{i=1}^N a_i \quad (26)$$

现目标函数为：

$$\min_{a_1, a_2} W(a_1, a_2) = \frac{1}{2} a_1^2 K_{11} + \frac{1}{2} a_2^2 K_{22} + a_1 a_2 y_1 y_2 K_{12} + c_1 y_1 a_1 + c_2 y_2 a_2 - a_1 - a_2 + const \quad (27)$$

其中 $c_i = \sum_{j=3}^N a_j y_j K_{ij}$ ， $const$ 为其他变量组成的常数，与优化目标无关，所以我们可以将目标函数写为：

$$\min_{a_1, a_2} W(a_1, a_2) = \frac{1}{2} a_1^2 K_{11} + \frac{1}{2} a_2^2 K_{22} + a_1 a_2 y_1 y_2 K_{12} + c_1 y_1 a_1 + c_2 y_2 a_2 - a_1 - a_2 \quad (28)$$

单变量约束

为了保证满足停机条件之一： $\sum_{i=1}^N a_i y_i = 0$ ，将 a_1 用 a_2 表示：

$$\begin{aligned}
a_1 y_1 + a_2 y_2 &= - \sum_{i=3}^N a_i y_i = z \\
a_1 &= y_1 (z - a_2 y_2)
\end{aligned} \quad (29)$$

将 $a_1 = z - a_2 y_2$ 带回到目标函数得到：

$$\min_{a_2} W(a_2) = \frac{1}{2} y_1^2 (z - a_2 y_2)^2 K_{11} + \frac{1}{2} a_2^2 K_{22} + y_1 (z - a_2 y_2) a_2 y_2 K_{12} + c_1 y_1^2 (z - a_2 y_2) + c_2 y_2 a_2 - y_1 (z - a_2 y_2) - a_2 \quad (30)$$

并且 $y_i^2 = 1$ ，可将上式化简得：

$$\min_{a_2} W(a_2) = \frac{1}{2} (z - a_2 y_2)^2 K_{11} + \frac{1}{2} a_2^2 K_{22} + (z - a_2 y_2) a_2 y_2 K_{12} + c_1 (z - a_2 y_2) + c_2 y_2 a_2 - y_1 (z - a_2 y_2) - a_2 \quad (31)$$

令其导数为0得到：

$$\begin{aligned}
\frac{dW}{da_2} &= -y_2 K_{11} (z - a_2 y_2) + K_{22} a_2 + K_{12} (y_2 (z - a_2 y_2) - y_2^2 a_2) - c_1 y_2 + c_2 y_2 + y_1 y_2 - 1 = 0 \\
&= (K_{11} - 2K_{12} + K_{22}) a_2 - y_2 (K_{11} z - 2K_{12} z + c_1 - c_2 - y_1 + y_2) = 0 \\
&= (K_{11} - 2K_{12} + K_{22}) a_2 = y_2 (K_{11} z - 2K_{12} z + c_1 - c_2 - y_1 + y_2)
\end{aligned}$$

由于是迭代算法，每一步需要更新 a_1, a_2 的值，所以上述的 z, c_i 为：

$$\begin{aligned} z &= a_1^{old} y_1 + a_2^{old} y_2 \\ c_i &= \sum_{j=3}^N a_j y_j K_{ij} = g(x_i^{old}) - b^{old} - a_1^{old} y_1 K_{1i} - a_2^{old} y_2 K_{2i} \end{aligned} \quad (32)$$

带回去得到：

$$\begin{aligned} (K_{11} - 2K_{12} + K_{22})a_2 &= y_2(K_{11}z - 2K_{12}z + c_1 - c_2 - y_1 + y_2) \\ (K_{11} - 2K_{12} + K_{22})a_2 &= y_2(K_{11}(a_1^{old}y_1 + a_2^{old}y_2) - 2K_{12}(a_1^{old}y_1 + a_2^{old}y_2) + (g(x_1^{old}) - b^{old} - a_1^{old}y_1K_{11} - a_2^{old}y_2K_{21}) - (g(x_2^{old}) - b^{old} - a_1^{old}y_1K_{12} - a_2^{old}y_2K_{22}) - y_1 + y_2) \\ (K_{11} - 2K_{12} + K_{22})a_2 &= (K_{11} - 2K_{12} + K_{22})a_2^{old} + y_2(g(x_1^{old}) - y_1) - y_2(g(x_2^{old}) - y_2) \end{aligned}$$

令 $\eta = K_{11} - 2K_{12} + K_{22}$ ，记 $E_i = g(x_i) - y_i$ 为第 i 个样本的差值。

所以有

$$a_2^{new} = a_2^{old} + y_2 \frac{E_1 - E_2}{\eta} \quad (33)$$

修剪

上述的 a_2^{new} 不一定满足 $0 \leq a_2^{new} \leq C$ 的条件，所以我们需要对其进行修剪：

等式约束：

$$a_1^{new} y_1 + a_2^{new} y_2 = a_1^{old} y_1 + a_2^{old} y_2 \quad (34)$$

可以反解出 a_1^{new} ：

$$a_1^{new} = y_1(a_1^{old} y_1 + a_2^{old} y_2 - a_2^{new} y_2) \quad (35)$$

不等式约束：

$$0 \leq a_1^{new} \leq C \quad (36)$$

- y_1, y_2 同号时：

$$\begin{aligned} a_1^{new} &= a_1^{old} - a_2^{old} + a_2^{new} \\ 0 \leq a_1^{new} \leq C &\Leftrightarrow a_2^{old} - a_1^{old} \leq a_2^{new} \leq a_2^{old} - a_1^{old} + C \end{aligned} \quad (37)$$

为了保证 $0 \leq a_2^{new} \leq C$ ，令：

$$\begin{aligned} L &= \max(0, a_2^{old} - a_1^{old}) \\ H &= \min(C, a_2^{old} - a_1^{old} + C) \end{aligned}$$

所以必须有：

$$L \leq a_2^{new} \leq H \quad (38)$$

当 a_2^{new} 越过边界时，取 a_2^{new} 为边界

- y_1, y_2 异号时：

$$\begin{aligned} a_1^{new} &= a_1^{old} + a_2^{old} - a_2^{new} \\ 0 \leq a_1^{new} \leq C &\Leftrightarrow a_2^{old} + a_1^{old} - C \leq a_2^{new} \leq a_2^{old} + a_1^{old} \end{aligned} \quad (39)$$

同理，令：

$$\begin{aligned} L &= \max(0, a_2^{old} + a_1^{old} - C) \\ H &= \min(C, a_2^{old} + a_1^{old}) \end{aligned}$$

所以必须有：

$$L \leq a_2^{new} \leq H \quad (40)$$

当 a_2^{new} 越过边界时，取 a_2^{new} 为边界

即：

$$a_2^{new} = \begin{cases} H, & \text{if } a_2^{new} \geq H \\ a_2^{new}, & \text{if } L < a_2^{new} < H \\ L, & \text{if } a_2^{new} \leq L \end{cases}$$

计算 a_1^{new}

根据计算出的 a_2^{new} 计算 a_1^{new}

$$a_1^{new} = y_1(a_1^{old}y_1 + a_2^{old}y_2 - a_2^{new}y_2) \quad (41)$$

选择变量的启发式方法

SMO在每一步都是选择两个变量进行优化，其中至少有一个违反KKT条件；根据Osuna的理论，每一步优化都会使目标函数值减小，因此可以保证算法的收敛性。

为加速收敛，SMO采用启发式的方法来选择每一步中优化的两个变量。有两个启发式规则：

- 一个是选择第一个变量的规则，构成外层循环；
- 另一个是选择第二个变量的规则，构成内层循环。

选择第一个优化变量

第一个优化变量选取训练样本中违反KKT条件最严重的样本点，并将其对应的变量作为第1个变量。具体地，检验样本点 (x_i, y_i) 是否满足：

$$\begin{aligned} a_i = 0 &\Leftrightarrow y_i g(x_i) \geq 1 \\ a_i = C &\Leftrightarrow y_i g(x_i) \leq 1 \\ 0 < a_i < C &\Leftrightarrow y_i g(x_i) = 1 \end{aligned}$$

其中 $g(x_i) = \sum_{j=1}^N a_j y_j K_{ij} + b$ 。有时为了加快收敛速度，会在 ξ 范围内进行的，即：

$$\begin{aligned} a_i = 0 &\Leftrightarrow y_i g(x_i) \geq 1 - \xi \\ a_i = C &\Leftrightarrow y_i g(x_i) \leq 1 + \xi \\ 0 < a_i < C &\Leftrightarrow 1 - \xi \leq y_i g(x_i) \leq 1 + \xi \end{aligned}$$

一般取 $\xi = 10^{-3} \sim 10^{-2}$ 。

在检验过程中，外层循环首先遍历所有满足条件 $0 < a_i < C$ 的样本点，即在间隔边界上的支持向量点，检验他们是否满足KKT条件，如果他们满足则遍历整个训练集，检验是否满足KKT条件。

选择第二个优化变量

选定第一个优化变量 a_1 之后，选择第二个优化变量 a_2 ，根据公式 $a_2^{new} = a_2^{old} + y_2 \frac{E_1 - E_2}{\eta}$ ，我们可以使优化的步长最大的样本点，也就是让 $y_2 \frac{E_1 - E_2}{\eta}$ 的绝对值最大，但是计算 η 需要计算核函数，为了简便，将 $|E_1 - E_2|$ 近似为步长，根据第一个优化变量可以确定 E_1 ，然后根据 E_1 的符号可以选择 E_2 使得 $|E_1 - E_2|$ 最大，具体地：

$$\begin{aligned} \text{if } E_1 < 0 : & \max(E_2) \\ \text{else : } & \min(E_2) \end{aligned} \quad (42)$$

当 E_1 小于0选取最大的作为 E_2 ，反之选最小的作为 E_2 。

需要注意的是，每个变量不一定仅仅被优化一次，可能跟随着其他变量被优化好多次。

E和b的更新

这位大佬写的很好，可以去借鉴一下，这里就不再过多叙述，只需要记住更新公式：

$$\begin{aligned}
b_1^{new} &= b^{old} - E_1^{old} - y_1 K_{11}(a_1^{new} - a_1^{old}) - y_2 K_{12}(a_2^{new} - a_2^{old}) \\
b_2^{new} &= b^{old} - E_2^{old} - y_1 K_{12}(a_1^{new} - a_1^{old}) - y_2 K_{22}(a_2^{new} - a_2^{old}) \\
E_i^{new} &= E_i^{old} + y_1 K_{1i}(a_1^{new} - a_1^{old}) + y_2 K_{2i}(a_2^{new} - a_2^{old}) + b^{new} - b^{old}
\end{aligned} \tag{43}$$

Python + Numpy实现

```

1  import numpy as np
2
3
4  class SVM():
5
6      def __init__(self,max_iter = 200,C = 1):
7          # 最多迭代数
8          self.max_iter = max_iter
9          # w
10         self.w = 0
11         # b
12         self.b = 0
13         # 松弛变量
14         self.C = C
15         # 精度(xi为其希腊字母)
16         self.xi = 0.01
17
18     def fit(self,X_train:np.ndarray,Y_train:np.ndarray):
19         # 样本个数和x的维度
20         self.m,self.n = X_train.shape
21         # X
22         self.X = X_train
23         # Y
24         self.Y = Y_train
25         # alpha
26         self.alpha = np.zeros(self.m)
27         # error_i = w·x_i+b-y_i
28         self.E = np.array([self.g(i)-self.Y[i] for i in range(self.m)])
29
30         for _iter in range(self.max_iter):
31
32             # 非边界alpha
33             non_bound_alpha = [i for i in range(self.m) if 0 <
self.alpha[i] < self.C]
34             # 边界alpha
35             bound_alpha = [i for i in range(self.m) if i not in
non_bound_alpha]
36
37             # 先非边界alpha后边界alpha
38             non_bound_alpha.extend(bound_alpha)
39
40             # 选取的两个变量的下标
41             i1 = i2 = None
42
43             for i in non_bound_alpha:
44                 # 满足KKT条件则跳出
45                 if self.KKT(i):
46                     continue
47
48                 # 第一个优化alpha下标

```

```

49         i1 = i
50         # 使第二个优化变量优化步长最大
51
52
53         if self.E[i] < 0:
54             i2 = self.E.argmax()
55         else:
56             i2 = self.E.argmin()
57         break
58
59         # 待优化变量
60         E1_old = self.E[i1]
61         E2_old = self.E[i2]
62         alpha1_old = self.alpha[i1]
63         alpha2_old = self.alpha[i2]
64         Y1 = self.Y[i1]
65         Y2 = self.Y[i2]
66
67         # eta(希腊字母) =  $\kappa_{11}-2\kappa_{12}+\kappa_{22}$ 
68         eta =
self.kernel(i1,i1)-2*self.kernel(i1,i2)+self.kernel(i2,i2)
69
70         # eta小于等于0则重新找优化变量
71         if eta <=0 :
72             continue
73
74         # 找到L,H
75         if Y1 == Y2:
76             L = max(0,alpha1_old+alpha2_old-self.C)
77             H = min(self.C,alpha1_old+alpha2_old)
78         else:
79             L = max(0,alpha2_old-alpha1_old)
80             H = min(self.C,alpha2_old-alpha1_old+self.C)
81
82         # 更新alpha2
83         alpha2_new_unclipped = alpha2_old + Y2*(E1_old-E2_old)/eta
84         alpha2_new = self.update_alpha2(alpha2_new_unclipped,L,H)
85
86         # 更新alpha1
87         alpha1_new = alpha1_old + Y1*Y2*(alpha2_old-alpha2_new)
88
89         # 更新b
90         b1_new = self.b-E1_old-Y1*self.kernel(i1,i1)*(alpha1_new-
alpha1_old)-Y2*self.kernel(i1,i2)*(alpha2_new-alpha2_old)
91         b2_new = self.b - E2_old-Y1*self.kernel(i1,i2)*(alpha1_new-
alpha1_old)-Y2*self.kernel(i2,i2)*(alpha2_new-alpha2_old)
92
93         if 0< alpha1_new <self.C:
94             b_new = b1_new
95         elif 0<alpha2_new <self.C:
96             b_new = b2_new
97         else:
98             b_new = (b1_new+b2_new)/2
99
100         # 更新E
101         self.E = self.E + Y1*(alpha1_new-
alpha1_old)*np.array([self.kernel(i1,j) for j in range(self.m)])+ \

```

```

102         Y2*(alpha2_new-alpha2_old)*np.array([self.kernel(i2,j) for
j in range(self.m)]) + np.array([b_new-self.b for _ in range(self.m)])
103
104         # 修改到原变量中
105         self.alpha[i1] = alpha1_new
106         self.alpha[i2] = alpha2_new
107         self.b = b_new
108
109         # 获取w
110         self.w = np.dot(self.alpha,self.Y.reshape(-1,1)*self.X)
111
112
113         # 更新alpha2
114         def update_alpha2(self, alpha2_new_unclipped, L, H):
115             # 越界则取边界
116             if alpha2_new_unclipped < L:
117                 return L
118             elif alpha2_new_unclipped > H:
119                 return H
120             # 没有越界则不变
121             else:
122                 return alpha2_new_unclipped
123
124         # KKT条件
125         def KKT(self,i):
126             if self.alpha[i] == 0:
127                 return self.Y[i]*self.g(i) >= 1-self.xi
128             elif self.alpha[i] == self.C:
129                 return self.Y[i]*self.g(i)<= 1+self.xi
130             else:
131                 return 1-self.xi<= self.Y[i]*self.g(i)<= 1+self.xi
132
133         # g(i) = w·x+b = SUM(a_j*y_j*k_ij)
134         def g(self,i):
135             sum = 0
136             for j in range(self.m):
137                 sum += self.alpha[j]*self.Y[j]*self.kernel(i,j)
138             return sum+self.b
139
140         # 核函数
141         def kernel(self,i,j):
142             return np.dot(self.X[i],self.X[j])
143
144
145         def predict(self,X):
146             return np.sign(np.dot(self.w,X)+self.b)
147
148         def score(self,X_test,Y_test):
149             right_num = 0
150             for i in range(len(X_test)):
151                 y = self.predict(X_test[i])
152                 if y == Y_test[i]:
153                     right_num+=1
154             return right_num/len(X_test)

```