

# 机器学习实验报告

## 支持向量机

朱天泽

(日期: 2022 年 4 月 17 日)

**摘要** 在《机器学习》第 6 章中, 我学习了支持向量机。此次实验, 我实现了使用核函数和软间隔的支持向量机, 在西瓜数据集 3.0 $\alpha$  上用多个核函数进行了训练, 并对不同核函数的分类效果进行了比较。

**关键词** SVM; 核函数; 分类

## 1 题目理解

题目要求: 在西瓜数据集 3.0 $\alpha$  上分别用线性核和高斯核训练一个 SVM。

针对不同的核函数  $\kappa$ , 其底层的数学原理都类似, 给一个 SVM 替换核函数也相对方便。因此本次实验计划不止用线性核和高斯核训练 SVM, 而是使用多个核函数分别训练一个 SVM, 并对不同核函数得到的 SVM 的分类效果做对比。

## 2 支持向量机原理阐述

### 2.1 间隔与支持向量

给定训练样本集  $D = \{(\mathbf{x}_1, y_1), (\mathbf{x}_2, y_2), \dots, (\mathbf{x}_m, y_m)\}$ ,  $y_i \in \{-1, 1\}$ , 分类学习最基本的想法是基于训练集  $D$  在样本空间中找到一个划分超平面, 将不同类别的样本分开。

在样本空间中, 划分超平面可通过如下线性方程来描述:

$$\mathbf{w}^T \mathbf{x} + b = 0 \quad (1)$$

其中  $\mathbf{w} = (w_1, w_2, \dots, w_d)^T$  为超平面的法向量,  $b$  为位移项, 决定了超平面与原点之间的距离。显然, 划分超平面可被法向量  $\mathbf{w}$  和位移  $b$  确定, 下面我们将超平面记为  $(\mathbf{w}, b)$ 。

对于任意一点  $\mathbf{x}_0 = (x_1^0, x_2^0, \dots, x_n^0)^T$ , 设其在超平面  $\mathbf{w}^T \mathbf{x} + b = 0$  上的投影点为  $\mathbf{x}_1 = (x_1^1, x_2^1, \dots, x_n^1)^T$ , 则  $\mathbf{w}^T \mathbf{x}_1 + b = 0$ , 设其到超平面的距离为  $r$ ; 由于向量  $\overrightarrow{\mathbf{x}_1 \mathbf{x}_0}$  与法向量  $\mathbf{w}$  平行, 因此

$$|\mathbf{w} \cdot \overrightarrow{\mathbf{x}_1 \mathbf{x}_0}| = \|\mathbf{w}\| \cdot \cos \pi \cdot \|\overrightarrow{\mathbf{x}_1 \mathbf{x}_0}\| = \|\mathbf{w}\| \cdot \|\overrightarrow{\mathbf{x}_1 \mathbf{x}_0}\| = \|\mathbf{w}\| \cdot r \quad (2)$$

又有

$$\begin{aligned} \mathbf{w} \cdot \overrightarrow{\mathbf{x}_1 \mathbf{x}_0} &= w_1 (x_1^0 - x_1^1) + w_2 (x_2^0 - x_2^1) + \dots + w_n (x_n^0 - x_n^1) \\ &= w_1 x_1^0 + w_2 x_2^0 + \dots + w_n x_n^0 - (w_1 x_1^1 + w_2 x_2^1 + \dots + w_n x_n^1) \\ &= \mathbf{w}^T \mathbf{x}_0 - \mathbf{w}^T \mathbf{x}_1 \\ &= \mathbf{w}^T \mathbf{x}_0 + b \end{aligned} \quad (3)$$

故  $|\mathbf{w}^T \mathbf{x}_0 + b| = \|\mathbf{w}\| \cdot r$ ，因此对于任意点  $\mathbf{x}$ ，其到超平面  $(\mathbf{w}, b)$  的距离为

$$r = \frac{|\mathbf{w}^T \mathbf{x} + b|}{\|\mathbf{w}\|} \quad (4)$$

将训练样本分开的划分超平面有很多，如图1。直观上看应该去找位于两类训练样本“正中间”的划分超平面，即距离超平面最近的正反例到超平面的距离“相当”，该划分超平面对训练样本局部扰动的容忍性最好。例如，由于训练集的局限性或噪声的因素，训练集外的样本可能比图1中的训练样本更接近两个类的分隔界，这将使许多划分超平面出现错误，而“正中间”的超平面受影响最小。换言之，这个划分超平面所产生的分类结果是最鲁棒的，对未见示例的泛化能力最强。

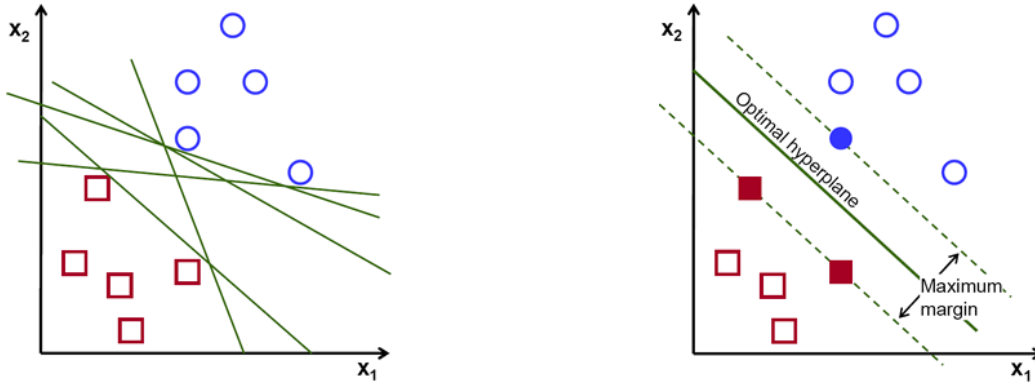


图 1: 有多个划分超平面将两类训练样本分开

对于给定的数据集  $D$  和超平面  $(\mathbf{w}, b)$ ，定义任意一个样本点  $(\mathbf{x}_i, y_i)$ ， $y_i \in \{-1, 1\}$ ， $i = 1, 2, \dots, m$  关于超平面的几何间隔为

$$\gamma_i = \frac{y_i (\mathbf{w}^T \mathbf{x}_i + b)}{\|\mathbf{w}\|} \quad (5)$$

当对样本  $(\mathbf{x}_i, y_i)$  正确分类时， $\gamma_i > 0$ ，几何间隔此时也等价于点到超平面的距离；而没有正确分类时， $\gamma_i < 0$ 。

对于给定的数据集  $D$  和超平面  $(\mathbf{w}, b)$ ，定义数据集  $D$  关于超平面的几何间隔为数据集  $D$  中所有样本点的几何间隔最小值

$$\gamma = \min_{i=1,2,\dots,m} \gamma_i \quad (6)$$

给定线性可分数据集  $D$ ，支持向量机模型希望求得使得数据集  $D$  的几何间隔  $\gamma$  达到最大的超平面，然后利用  $\text{sign}$  函数实现分类功能

$$y = \text{sign}(\mathbf{w}^T \mathbf{x} + b) = \begin{cases} 1, & \mathbf{w}^T \mathbf{x} + b > 0 \\ -1, & \mathbf{w}^T \mathbf{x} + b < 0 \end{cases} \quad (7)$$

易知：

- 当超平面没有正确划分样本时，必存在  $\gamma_i < 0$ ，故  $\gamma < 0$ ；
- 当超平面正确划分样本时，有  $\gamma \geq 0$ ，且越靠近样本间隔中央  $\gamma$  越大

因此，使得几何间隔  $\gamma$  最大的超平面就是位于两类训练样本“正中间”的划分超平面。

给定线性可分数据集  $D$ ，设  $D$  中几何间隔最小的样本为  $(\mathbf{x}_{\min}, y_{\min})$ ，称其为“支持向量”，如图2所示。

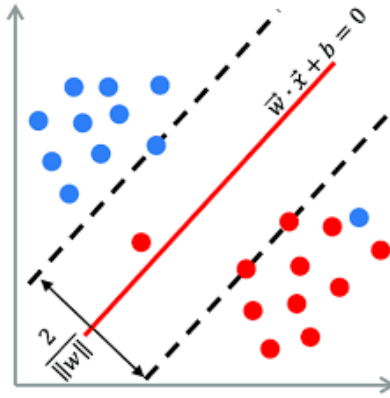


图 2: 支持向量与几何间隔

SVM 寻找最优超平面的过程可以转化为如下带约束条件的优化问题

$$\begin{aligned} \max \quad & \gamma \\ \text{s.t.} \quad & \gamma_i \geq \gamma, \quad i = 1, 2, \dots, m \end{aligned} \quad (8)$$

即

$$\begin{aligned} \max_{\mathbf{w}, b} \quad & \frac{y_{\min}(\mathbf{w}^T \mathbf{x}_{\min} + b)}{\|\mathbf{w}\|} \\ \text{s.t.} \quad & y_i(\mathbf{w}^T \mathbf{x}_i + b) \geq y_{\min}(\mathbf{w}^T \mathbf{x}_{\min} + b), \quad i = 1, 2, \dots, m \end{aligned} \quad (9)$$

假设该问题的最优解为  $(\mathbf{w}^*, b^*)$ , 那么  $(k\mathbf{w}^*, kb^*)$ ,  $k \in \mathbb{R}^+$  也是最优解, 且超平面不变, 因此需要对  $(\mathbf{w}, b)$  做一定限制才能使得上述优化问题有唯一解。那么不妨令  $y_{\min}(\mathbf{w}^T \mathbf{x}_{\min} + b) = 1$ , 对于特定的  $(\mathbf{x}_{\min}, y_{\min})$  来说, 能使得  $y_{\min}(\mathbf{w}^T \mathbf{x}_{\min} + b) = 1$  的  $k$  有且仅有一个。因此, 上述优化问题进一步转化为

$$\begin{aligned} \max_{\mathbf{w}, b} \quad & \frac{1}{\|\mathbf{w}\|} \\ \text{s.t.} \quad & y_i(\mathbf{w}^T \mathbf{x}_i + b) \geq 1, \quad i = 1, 2, \dots, m \end{aligned} \quad (10)$$

上述问题等价于

$$\begin{aligned} \min_{\mathbf{w}, b} \quad & \frac{1}{2} \|\mathbf{w}\|^2 \\ \text{s.t.} \quad & 1 - y_i(\mathbf{w}^T \mathbf{x}_i + b) \leq 0, \quad i = 1, 2, \dots, m \end{aligned} \quad (11)$$

这就是 SVM 的基本型。

## 2.2 对偶问题

我们称式(11)为主问题, 对式(11)使用拉格朗日乘子法, 其拉格朗日函数为

$$L(\mathbf{w}, b, \boldsymbol{\alpha}) = \frac{1}{2} \|\mathbf{w}\|^2 + \sum_{i=1}^m \alpha_i (1 - y_i(\mathbf{w}^T \mathbf{x}_i + b)) \quad (12)$$

其中  $\boldsymbol{\alpha} = (\alpha_1, \alpha_2, \dots, \alpha_m)^T$ 。令  $L(\mathbf{w}, b, \boldsymbol{\alpha})$  对  $\mathbf{w}$  和  $b$  的偏导数都为零, 可得

$$\mathbf{w} = \sum_{i=1}^m \alpha_i y_i \mathbf{x}_i \quad (13)$$

$$\sum_{i=1}^m \alpha_i y_i = 0 \quad (14)$$

将式(13)代入(12)

$$\begin{aligned} \inf_{\mathbf{w}, b} L(\mathbf{w}, b, \boldsymbol{\alpha}) &= \frac{1}{2} \mathbf{w}^T \mathbf{w} + \sum_{i=1}^m \alpha_i - \sum_{i=1}^m \alpha_i y_i \mathbf{w}^T \mathbf{x}_i - \sum_{i=1}^m \alpha_i y_i b \\ &= \frac{1}{2} \mathbf{w}^T \sum_{i=1}^m \alpha_i y_i \mathbf{x}_i - \mathbf{w}^T \sum_{i=1}^m \alpha_i y_i \mathbf{x}_i + \sum_{i=1}^m \alpha_i - b \sum_{i=1}^m \alpha_i y_i \\ &= -\frac{1}{2} \mathbf{w}^T \sum_{i=1}^m \alpha_i y_i \mathbf{x}_i + \sum_{i=1}^m \alpha_i - b \sum_{i=1}^m \alpha_i y_i \end{aligned} \quad (15)$$

再将(14)代入(15), 于是有

$$\begin{aligned} \inf_{\mathbf{w}, b} L(\mathbf{w}, b, \boldsymbol{\alpha}) &= -\frac{1}{2} \mathbf{w}^T \sum_{i=1}^m \alpha_i y_i \mathbf{x}_i + \sum_{i=1}^m \alpha_i \\ &= -\frac{1}{2} \left( \sum_{i=1}^m \alpha_i y_i \mathbf{x}_i \right)^T \left( \sum_{i=1}^m \alpha_i y_i \mathbf{x}_i \right) + \sum_{i=1}^m \alpha_i \\ &= -\frac{1}{2} \sum_{i=1}^m \alpha_i y_i \mathbf{x}_i^T \sum_{i=1}^m \alpha_i y_i \mathbf{x}_i + \sum_{i=1}^m \alpha_i \\ &= \sum_{i=1}^m \alpha_i - \frac{1}{2} \sum_{i=1}^m \sum_{j=1}^m \alpha_i \alpha_j y_i y_j \mathbf{x}_i^T \mathbf{x}_j \end{aligned} \quad (16)$$

最终得到式(11)的对偶问题

$$\begin{aligned} \max_{\boldsymbol{\alpha}} \quad & \sum_{i=1}^m \alpha_i - \frac{1}{2} \sum_{i=1}^m \sum_{j=1}^m \alpha_i \alpha_j y_i y_j \mathbf{x}_i^T \mathbf{x}_j \\ \text{s.t.} \quad & \sum_{i=1}^m \alpha_i y_i = 0 \\ & \alpha_i \geq 0, \quad i = 1, 2, \dots, m \end{aligned} \quad (17)$$

注意到主问题, 即式(11), 是凸优化问题, 且除支持向量外, 存在样本点满足  $1 - y_i(\mathbf{w}^T \mathbf{x}_i + b) < 0$ ; 因此强对偶性成立, 对偶问题的解就是原问题的解, 且必须满足 KKT 条件

$$\begin{cases} \alpha_i \geq 0 \\ y_i f(\mathbf{x}_i) - 1 \geq 0 \\ \alpha_i (y_i f(\mathbf{x}_i) - 1) = 0 \end{cases} \quad (18)$$

解得  $\boldsymbol{\alpha}$  后, 求出  $\mathbf{w}$  和  $b$  即可得到模型

$$\begin{aligned} f(\mathbf{x}) &= \mathbf{w}^T \mathbf{x} + b \\ &= \sum_{i=1}^m \alpha_i y_i \mathbf{x}_i^T \mathbf{x} + b \end{aligned} \quad (19)$$

## 2.3 核函数

在前面的讨论中，我们假设训练样本是线性可分的，即存在一个划分超平面能将训练样本正确分类。然而在现实任务中，原始样本空间内也许并不存在一个能正确划分两类样本的超平面。对这样的问题，可将样本从原始空间映射到一个更高维的特征空间，使得样本在这个特征空间内线性可分。

令  $\phi(\mathbf{x})$  表示将  $\mathbf{x}$  映射后的特征向量，于是，在特征空间中划分超平面所对应的模型可表示为

$$f(\mathbf{x}) = \mathbf{w}^T \phi(\mathbf{x}) + b \quad (20)$$

其中  $\mathbf{w}$  和  $b$  是模型参数。类似式(11)，有

$$\begin{aligned} \min_{\mathbf{w}, b} \quad & \frac{1}{2} \|\mathbf{w}\|^2 \\ \text{s.t.} \quad & 1 - y_i (\mathbf{w}^T \phi(\mathbf{x}_i) + b) \leq 0, \quad i = 1, 2, \dots, m \end{aligned} \quad (21)$$

其对偶问题是

$$\begin{aligned} \max_{\alpha} \quad & \sum_{i=1}^m \alpha_i - \frac{1}{2} \sum_{i=1}^m \sum_{j=1}^m \alpha_i \alpha_j y_i y_j \phi(\mathbf{x}_i)^T \phi(\mathbf{x}_j) \\ \text{s.t.} \quad & \sum_{i=1}^m \alpha_i y_i = 0 \\ & \alpha_i \geq 0, \quad i = 1, 2, \dots, m \end{aligned} \quad (22)$$

求解式(22)涉及到计算  $\phi(\mathbf{x}_i)^T \phi(\mathbf{x}_j)$ ，这是样本  $\mathbf{x}_i$  与  $\mathbf{x}_j$  映射到特征空间之后的内积。由于特征空间维数可能很高，甚至可能是无穷维，因此直接计算  $\phi(\mathbf{x}_i)^T \phi(\mathbf{x}_j)$  通常是困难的。为了避免这样的障碍，可以使用核技巧，定义核函数

$$\kappa(\mathbf{x}_i, \mathbf{x}_j) = \langle \phi(\mathbf{x}_i), \phi(\mathbf{x}_j) \rangle = \phi(\mathbf{x}_i)^T \phi(\mathbf{x}_j) \quad (23)$$

即  $\mathbf{x}_i$  与  $\mathbf{x}_j$  在特征空间中的内积等于它们在原始样本空间中通过函数  $\kappa(\cdot, \cdot)$  计算的结果。有了核函数，我们就不必直接去计算高维甚至无穷维特征空间中的内积，于是式(22)可重写为

$$\begin{aligned} \max_{\alpha} \quad & \sum_{i=1}^m \alpha_i - \frac{1}{2} \sum_{i=1}^m \sum_{j=1}^m \alpha_i \alpha_j y_i y_j \kappa(\mathbf{x}_i, \mathbf{x}_j) \\ \text{s.t.} \quad & \sum_{i=1}^m \alpha_i y_i = 0, \\ & \alpha_i \geq 0, \quad i = 1, 2, \dots, m \end{aligned} \quad (24)$$

求解后即可得到

$$\begin{aligned} f(\mathbf{x}) &= \mathbf{w}^T \phi(\mathbf{x}) + b \\ &= \sum_{i=1}^m \alpha_i y_i \phi(\mathbf{x}_i)^T \phi(\mathbf{x}) + b \\ &= \sum_{i=1}^m \alpha_i y_i \kappa(\mathbf{x}, \mathbf{x}_i) + b \end{aligned} \quad (25)$$

**定理** 令  $\chi$  为输入空间， $\kappa(\cdot, \cdot)$  是定义在  $\chi \times \chi$  上的对称核函数，则  $\kappa$  是核函数当且仅当对于任意数据  $D = \mathbf{x}_{i=1}^m$ ，“核矩阵”  $\mathbf{K}$  总是半正定的：

$$\mathbf{K} = \begin{bmatrix} \kappa(\mathbf{x}_1, \mathbf{x}_1) & \cdots & \kappa(\mathbf{x}_1, \mathbf{x}_j) & \cdots & \kappa(\mathbf{x}_1, \mathbf{x}_m) \\ \vdots & \ddots & \vdots & \ddots & \vdots \\ \kappa(\mathbf{x}_i, \mathbf{x}_1) & \cdots & \kappa(\mathbf{x}_i, \mathbf{x}_j) & \cdots & \kappa(\mathbf{x}_i, \mathbf{x}_m) \\ \vdots & \ddots & \vdots & \ddots & \vdots \\ \kappa(\mathbf{x}_m, \mathbf{x}_1) & \cdots & \kappa(\mathbf{x}_m, \mathbf{x}_j) & \cdots & \kappa(\mathbf{x}_m, \mathbf{x}_m) \end{bmatrix}$$

在现实中，我们通常不知道  $\phi(\cdot)$  的具体形式，但是根据上述定理，我们不必知道  $\phi(\cdot)$  的具体形式，就能得到核函数；事实上，对于一个半正定核矩阵，总能找到一个与之对应的映射  $\phi$ 。我们希望样本在特征空间内线性可分，因此特征空间的好坏对支持向量机的性能至关重要。需注意的是，在不知道特征映射的形式时，我们并不知道什么样的核函数是合适的，而核函数也仅是隐式地定义了这个特征空间。于是，“核函数选择”成为 SVM 的最大变数。若核函数选择不合适，则意味着将样本映射到了一个不合适的特征空间，很可能导致性能不佳。

常见的核函数如表1所示。

名称	表达式	参数说明
线性核	$\kappa(\mathbf{x}_i, \mathbf{x}_j) = \mathbf{x}_i^T \mathbf{x}_j$	
多项式核	$\kappa(\mathbf{x}_i, \mathbf{x}_j) = (\mathbf{x}_i^T \mathbf{x}_j)^d$	$d \geq 1$ 为多项式的次数
高斯核	$\kappa(\mathbf{x}_i, \mathbf{x}_j) = e^{-\frac{\ \mathbf{x}_i - \mathbf{x}_j\ ^2}{2\sigma^2}}$	$\sigma > 0$ 为高斯核的带宽
RBF 核	$\kappa(\mathbf{x}_i, \mathbf{x}_j) = e^{-\gamma\ \mathbf{x}_i - \mathbf{x}_j\ ^2}$	$0 < \gamma < 1$
拉普拉斯核	$\kappa(\mathbf{x}_i, \mathbf{x}_j) = e^{-\frac{\ \mathbf{x}_i - \mathbf{x}_j\ }{\sigma}}$	$\sigma > 0$
Sigmoid 核	$\kappa(\mathbf{x}_i, \mathbf{x}_j) = \tanh(\beta \mathbf{x}_i^T \mathbf{x}_j + \theta)$	$\tanh$ 为双曲正切函数, $\beta > 0, \theta < 0$

表 1: 常用核函数

## 2.4 软间隔

在前面的讨论中，我们一直假定训练样本在样本空间或特征空间中是线性可分的，即存在一个超平面能将不同类的样本完全划分开。然而，在现实任务中往往很难确定合适的核函数使得训练样本在特征空间中线性可分。

缓解该问题的一个办法是允许支持向量机在一些样本上出错。为此，要引入“软间隔”的概念，如图3所示。

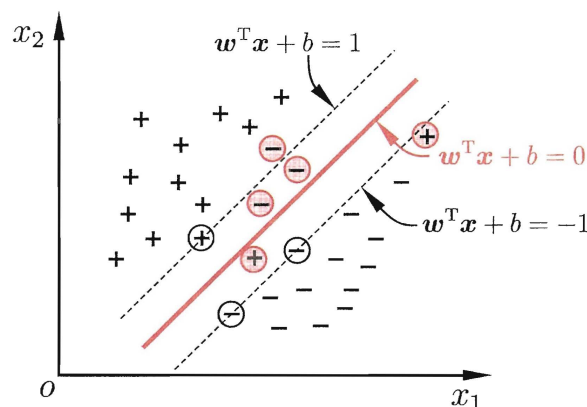


图 3: 软间隔示意图，红色圈出了一些不满足约束的样本。

以线性 SVM 为例，其形式要求所有样本都划分正确，即  $y_i (\mathbf{w}^T \mathbf{x}_i + b) \geq 1$ ，软间隔允许某些样本不满足这一约束。

当然，在最大化间隔的同时，不满足约束的样本应尽可能少。于是，优化目标可写为

$$\min_{\mathbf{w}, b} \frac{1}{2} \|\mathbf{w}\|^2 + C \sum_{i=1}^m \ell_{0/1} (y_i (\mathbf{w}^T \mathbf{x}_i + b) - 1) \quad (26)$$

其中  $C > 0$  是一个常数， $\ell_{0/1}$  是“0/1 损失函数”。 $\ell_{0/1}$  非凸、非连续，我们不妨使用 hinge 损失代替：

$$\ell_{\text{hinge}}(z) = \max(0, 1 - z) \quad (27)$$

则式(26)可写成

$$\min_{\mathbf{w}, b} \frac{1}{2} \|\mathbf{w}\|^2 + C \sum_{i=1}^m \max(0, 1 - y_i (\mathbf{w}^T \mathbf{x}_i + b)) \quad (28)$$

引入“松弛变量”  $\xi_i \geq 0$ ，可将式(28)重写为

$$\begin{aligned} \min_{\mathbf{w}, b, \xi_i} \quad & \frac{1}{2} \|\mathbf{w}\|^2 + C \sum_{i=1}^m \xi_i \\ \text{s.t.} \quad & y_i (\mathbf{w}^T \mathbf{x}_i + b) \geq 1 - \xi_i \\ & \xi_i \geq 0, \quad i = 1, 2, \dots, m \end{aligned} \quad (29)$$

这就是常用的软间隔 SVM。

式(29)与式(11)类似，仍是一个二次规划问题；于是，类似式(12)，通过拉格朗日乘子法可得到式(29)的拉格朗日函数

$$L(\mathbf{w}, b, \boldsymbol{\alpha}, \boldsymbol{\xi}, \boldsymbol{\mu}) = \frac{1}{2} \|\mathbf{w}\|^2 + C \sum_{i=1}^m \xi_i + \sum_{i=1}^m \alpha_i [1 - \xi_i - y_i (\mathbf{w}^T \mathbf{x}_i + b)] - \sum_{i=1}^m \mu_i \xi_i \quad (30)$$

其中  $\alpha_i \geq 0$ ， $\mu_i \geq 0$  是拉格朗日因子。

令  $L(\mathbf{w}, b, \boldsymbol{\alpha}, \boldsymbol{\xi}, \boldsymbol{\mu})$  对  $\mathbf{w}, b, \xi_i$  的偏导为零可得

$$\mathbf{w} = \sum_{i=1}^m \alpha_i y_i \mathbf{x}_i \quad (31)$$

$$\sum_{i=1}^m \alpha_i y_i = 0 \quad (32)$$

$$C = \alpha_i + \mu_i \quad (33)$$

将式(31)-(33)代入式(30)即可得到式(29)的对偶问题

$$\begin{aligned} \max_{\boldsymbol{\alpha}} \quad & \sum_{i=1}^m \alpha_i - \frac{1}{2} \sum_{i=1}^m \sum_{j=1}^m \alpha_i \alpha_j y_i y_j \mathbf{x}_i^T \mathbf{x}_j \\ \text{s.t.} \quad & \sum_{i=1}^m \alpha_i y_i = 0, \\ & 0 \leq \alpha_i \leq C, \quad i = 1, 2, \dots, m \end{aligned} \quad (34)$$

在线性的软间隔 SVM 的基础上，我们引入核函数  $\kappa$ ，对偶问题为

$$\begin{aligned}
& \max_{\alpha} \sum_{i=1}^m \alpha_i - \frac{1}{2} \sum_{i=1}^m \sum_{j=1}^m \alpha_i \alpha_j y_i y_j \kappa(\mathbf{x}_i, \mathbf{x}_j) \\
& \text{s.t.} \quad \sum_{i=1}^m \alpha_i y_i = 0, \\
& \quad 0 \leq \alpha_i \leq C, \quad i = 1, 2, \dots, m
\end{aligned} \tag{35}$$

类似式(18), 软间隔 SVM 的 KKT 条件为

$$\begin{cases} \alpha_i \geq 0, & \mu_i \geq 0 \\ y_i f(\mathbf{x}_i) - 1 + \xi_i \geq 0 \\ \alpha_i (y_i f(\mathbf{x}_i) - 1 + \xi_i) = 0 \\ \xi_i \geq 0, & \mu_i \xi_i = 0 \end{cases} \tag{36}$$

### 3 算法设计思路

#### 3.1 对偶问题转化

将式(35)表示的对偶问题转化为凸优化问题

$$\begin{aligned}
& \min_{\alpha} \frac{1}{2} \sum_{i=1}^m \sum_{j=1}^m \alpha_i \alpha_j y_i y_j \kappa(\mathbf{x}_i, \mathbf{x}_j) - \sum_{i=1}^m \alpha_i \\
& \text{s.t.} \quad \sum_{i=1}^m \alpha_i y_i = 0, \\
& \quad 0 \leq \alpha_i \leq C, \quad i = 1, 2, \dots, m
\end{aligned} \tag{37}$$

根据式(36)表示的 KKT 条件和式(33)表示的约束条件

$$\begin{aligned}
\alpha_i = 0 & \Rightarrow y_i f(\mathbf{x}_i) - 1 \geq 0 \\
0 < \alpha_i < C & \Rightarrow y_i f(\mathbf{x}_i) - 1 = 0 \\
\alpha_i = C & \Rightarrow y_i f(\mathbf{x}_i) - 1 \leq 0
\end{aligned} \tag{38}$$

考虑到计算机运算时的误差, 我们引入计算误差项  $\varepsilon \geq 0$ , 并记  $E_i = f(\mathbf{x}_i) - y_i$ , 得到

$$\begin{aligned}
\alpha_i = 0 & \Rightarrow y_i E_i \geq -\varepsilon \\
0 < \alpha_i < C & \Rightarrow |y_i E_i| < \varepsilon \\
\alpha_i = C & \Rightarrow y_i E_i \leq \varepsilon
\end{aligned} \tag{39}$$

进一步归纳, 得到  $\alpha_i$  的约束条件

$$\begin{aligned}
\alpha_i > 0 & \Rightarrow y_i E_i \leq \varepsilon \\
\alpha_i < C & \Rightarrow y_i E_i \geq -\varepsilon
\end{aligned} \tag{40}$$

#### 3.2 SMO 算法

##### 3.2.1 基本思想

SMO 的基本思路是先固定  $\alpha_i$  之外的所有参数, 然后求  $\alpha_i$  上的极值。由于存在约束  $\sum_{i=1}^m \alpha_i y_i = 0$ , 若固定  $\alpha_i$  之外的其它变量, 则  $\alpha_i$  可由其它变量导出。于是, SMO 算法每次选择两个变量  $\alpha_i$  和  $\alpha_j$ , 并固定其



它参数。这样，在参数初始化后，SMO 不断执行如下步骤直至收敛：

- 选取一对需要更新的  $\alpha_i$  和  $\alpha_j$ ；
- 固定  $\alpha_i$  和  $\alpha_j$  以外的参数，求解式(37)获得更新后的  $\alpha_i$  和  $\alpha_j$ 。

### 3.2.2 单变量无约束迭代

下面来推导更新  $\alpha_i$  和  $\alpha_j$  的迭代式。由于  $\alpha$  中所有维度都是等价的，我们不妨简化问题，选取  $\alpha_1$  和  $\alpha_2$ ，推导更新  $\alpha_i$  和  $\alpha_j$  的迭代式。

令

$$\Psi(\alpha) = \frac{1}{2} \sum_{i=1}^m \sum_{j=1}^m \alpha_i \alpha_j y_i y_j \kappa(\mathbf{x}_i, \mathbf{x}_j) - \sum_{i=1}^m \alpha_i \quad (41)$$

我们选取  $\alpha_1$  和  $\alpha_2$ ，固定  $\alpha$  中其它维度，省略  $\Psi(\alpha)$  中的常数项，得到

$$\begin{aligned} \Psi(\alpha_1, \alpha_2) &= \frac{1}{2} \alpha_1^2 \kappa(\mathbf{x}_1, \mathbf{x}_1) + \frac{1}{2} \alpha_2^2 \kappa(\mathbf{x}_2, \mathbf{x}_2) \\ &\quad + y_1 y_2 \alpha_1 \alpha_2 \kappa(\mathbf{x}_1, \mathbf{x}_2) + \nu_1 y_1 \alpha_1 + \nu_2 y_2 \alpha_2 - \alpha_1 - \alpha_2 \end{aligned} \quad (42)$$

其中

$$\begin{aligned} \nu_i &= \sum_{j=3}^m \alpha_j y_j \kappa(\mathbf{x}_i, \mathbf{x}_j) \\ &= f(\mathbf{x}_i) - y_1 \alpha_1 \kappa(\mathbf{x}_i, \mathbf{x}_1) - y_2 \alpha_2 \kappa(\mathbf{x}_i, \mathbf{x}_2) - b, \quad i = 1, 2 \end{aligned} \quad (43)$$

根据约束条件式(32)，有

$$\alpha_1 y_1 + \alpha_2 y_2 = - \sum_{i=3}^m \alpha_i y_i = \zeta \quad (44)$$

其中  $\zeta$  为常数。针对  $\alpha_1$  和  $\alpha_2$ ，假设更新之前为  $\alpha_1^{\text{old}}$  和  $\alpha_2^{\text{old}}$ ，更新之后为  $\alpha_1^{\text{new}}$  和  $\alpha_2^{\text{new}}$ ；由于  $\alpha$  中其余维度固定，因此必须满足

$$\alpha_1^{\text{old}} y_1 + \alpha_2^{\text{old}} y_2 = \alpha_1^{\text{new}} y_1 + \alpha_2^{\text{new}} y_2 = \zeta \quad (45)$$

两个变量不好同时求解，可以先求  $\alpha_2^{\text{new}}$ ，然后根据式(45)表示  $\alpha_1^{\text{new}}$ 。

由式(44)得

$$\alpha_1 = (\zeta - \alpha_2 y_2) y_1 \quad (46)$$

将式(46)代入式(41)，得

$$\begin{aligned} \Psi(\alpha_2) &= \frac{1}{2} (\zeta - \alpha_2 y_2)^2 \kappa(\mathbf{x}_1, \mathbf{x}_1) + \frac{1}{2} \alpha_2^2 \kappa(\mathbf{x}_2, \mathbf{x}_2) + \alpha_2 y_2 (\zeta - \alpha_2 y_2) \kappa(\mathbf{x}_1, \mathbf{x}_2) \\ &\quad + \nu_1 (\zeta - \alpha_2 y_2) + \nu_2 y_2 \alpha_2 - (\zeta - \alpha_2 y_2) y_1 - \alpha_2 \end{aligned} \quad (47)$$

令  $\frac{d\Psi(\alpha_2)}{d\alpha_2} = 0$ ，得

$$\alpha_2 = \frac{y_2 [y_2 - y_1 + \nu_1 - \nu_2 + \zeta (\kappa(\mathbf{x}_1, \mathbf{x}_1) - \kappa(\mathbf{x}_1, \mathbf{x}_2))]}{\kappa(\mathbf{x}_1, \mathbf{x}_1) + \kappa(\mathbf{x}_2, \mathbf{x}_2) - 2\kappa(\mathbf{x}_1, \mathbf{x}_2)} \quad (48)$$

将式(43)代入式(48)，再结合式(44)，得未修剪的  $\alpha_2^{\text{new}}$  为

$$\alpha_2^{\text{unclipped}} = \alpha_2^{\text{old}} + \frac{y_2(E_1 - E_2)}{\kappa(\mathbf{x}_1, \mathbf{x}_1) + \kappa(\mathbf{x}_2, \mathbf{x}_2) - 2\kappa(\mathbf{x}_1, \mathbf{x}_2)} \quad (49)$$

### 3.2.3 修剪 $\alpha_2^{\text{unclipped}}$

$\alpha_i$  必须满足不等式约束

$$0 \leq \alpha_i \leq C \quad (50)$$

式(49)未考虑到不等式约束，是未修剪的值。下面考虑不等式约束，由式(45)，有

$$\alpha_1^{\text{new}} = \alpha_1^{\text{old}} + y_1 y_2 (\alpha_2^{\text{old}} - \alpha_2^{\text{new}}) \quad (51)$$

当  $y_1 = y_2$ ，由(51)有

$$\alpha_1^{\text{new}} = \alpha_1^{\text{old}} + \alpha_2^{\text{old}} - \alpha_2^{\text{new}} \quad (52)$$

故

$$0 \leq \alpha_1^{\text{new}} \leq C \Leftrightarrow \alpha_1^{\text{old}} + \alpha_2^{\text{old}} - C \leq \alpha_2^{\text{new}} \leq \alpha_1^{\text{old}} + \alpha_2^{\text{old}} \quad (53)$$

又  $0 \leq \alpha_i \leq C$ ，所以

$$\max(0, \alpha_1^{\text{old}} + \alpha_2^{\text{old}} - C) \leq \alpha_2^{\text{new}} \leq \min(C, \alpha_1^{\text{old}} + \alpha_2^{\text{old}}) \quad (54)$$

同理，当  $y_1 \neq y_2$  时

$$\max(0, \alpha_2^{\text{old}} - \alpha_1^{\text{old}}) \leq \alpha_2^{\text{new}} \leq \min(C, \alpha_2^{\text{old}} - \alpha_1^{\text{old}} + C) \quad (55)$$

设  $[L, H]$  为考虑不等式约束后  $\alpha_2$  的可行域，整理得到

$$[L, H] = \begin{cases} [\max(0, \alpha_1^{\text{old}} + \alpha_2^{\text{old}} - C), \min(C, \alpha_1^{\text{old}} + \alpha_2^{\text{old}})], & y_1 = y_2 \\ [\max(0, \alpha_2^{\text{old}} - \alpha_1^{\text{old}}), \min(C, \alpha_2^{\text{old}} - \alpha_1^{\text{old}} + C)], & y_1 \neq y_2 \end{cases} \quad (56)$$

对  $\alpha_2^{\text{unclipped}}$  的修剪方法 clip 为

$$\alpha_2^{\text{new}} = \text{clip}(\alpha_2^{\text{unclipped}}) = \begin{cases} L, & \alpha_2^{\text{unclipped}} < L \\ \alpha_2^{\text{unclipped}}, & L \leq \alpha_2^{\text{unclipped}} \leq H \\ H, & \alpha_2^{\text{unclipped}} > H \end{cases} \quad (57)$$

### 3.2.4 计算 $\alpha_1^{\text{new}}$

由于式(56)是从式(51)开始，根据不等式约束  $0 \leq \alpha_i \leq C$  得到，故由式(51)得到的  $\alpha_1^{\text{new}}$  必然满足不等式约束， $\alpha_1$  的迭代式就为式(51)。

### 3.2.5 计算 $b^{\text{new}}$

若  $0 < \alpha_1^{\text{new}} < C$ , 由式(25), 有

$$\begin{aligned} b^{\text{new}} = b_1^{\text{new}} &= y_1 - \sum_{i=3}^m \alpha_i y_i \kappa(\mathbf{x}_1, \mathbf{x}_i) - \alpha_1^{\text{new}} y_1 \kappa(\mathbf{x}_1, \mathbf{x}_1) - \alpha_2^{\text{new}} y_2 \\ &= (\alpha_1^{\text{old}} - \alpha_1^{\text{new}}) y_1 \kappa(\mathbf{x}_1, \mathbf{x}_1) + (\alpha_2^{\text{old}} - \alpha_2^{\text{new}}) y_1 \kappa(\mathbf{x}_1, \mathbf{x}_2) - E_1^{\text{old}} + b^{\text{old}} \end{aligned} \quad (58)$$

同理, 若  $0 < \alpha_2^{\text{new}} < C$

$$b^{\text{new}} = b_2^{\text{new}} = (\alpha_1^{\text{old}} - \alpha_1^{\text{new}}) y_1 \kappa(\mathbf{x}_1, \mathbf{x}_2) + (\alpha_2^{\text{old}} - \alpha_2^{\text{new}}) y_1 \kappa(\mathbf{x}_2, \mathbf{x}_2) - E_2^{\text{old}} + b^{\text{old}} \quad (59)$$

若同时满足  $0 < \alpha_1^{\text{new}} < C$  和  $0 < \alpha_2^{\text{new}} < C$ , 则显然有  $b^{\text{new}} = b_1^{\text{new}} = b_2^{\text{new}}$ ; 若  $\alpha_1^{\text{new}}, \alpha_2^{\text{new}} \in \{0, C\}$  且  $L \neq H$ , 则  $b_1^{\text{new}}, b_2^{\text{new}}$  之间的值都能使两个样本满足 KKT 条件, 取  $b^{\text{new}} = \frac{b_1^{\text{new}} + b_2^{\text{new}}}{2}$ ; 若  $\alpha_1^{\text{new}}, \alpha_2^{\text{new}} \in \{0, C\}$  且  $L = H$ , 则不更新  $b$ 。

### 3.2.6 启发式选择

根据前面的推导, 每一步选择两个维度  $\alpha_i$  和  $\alpha_j$  进行优化

$$\begin{aligned} \alpha_i^{\text{new}} &= \alpha_i^{\text{old}} + y_1 y_2 (\alpha_j^{\text{old}} - \alpha_j^{\text{new}}) \\ \alpha_j^{\text{new}} &= \text{clip} \left( \alpha_j^{\text{old}} + \frac{y_j (E_i - E_j)}{\kappa(\mathbf{x}_i, \mathbf{x}_i) + \kappa(\mathbf{x}_j, \mathbf{x}_j) - 2\kappa(\mathbf{x}_i, \mathbf{x}_j)} \right) \end{aligned} \quad (60)$$

注意到只需选取的  $\alpha_i$  和  $\alpha_j$  中有一个不满足式(40)的约束, 目标函数就会在迭代后减小。所以 SMO 算法先选取违背式(40)的维度, 得到  $\alpha_i$ 。选择  $\alpha_j$  的启发规则是选择样本  $\kappa(\mathbf{x}_j, y_j)$  来使得优化的步长最大化, 也就是使式(57)中  $\frac{y_j (E_i - E_j)}{\kappa(\mathbf{x}_i, \mathbf{x}_i) + \kappa(\mathbf{x}_j, \mathbf{x}_j) - 2\kappa(\mathbf{x}_i, \mathbf{x}_j)}$  最大化。由于计算核函数比较耗时间, 所以 SMO 使用简单的近似, 用  $|E_i - E_j|$  来近似步长; 也就是说, 选择样本  $(\mathbf{x}_j, y_j)$  使  $|E_i - E_j|$  最大,  $\alpha_j$  的更新步长就能尽可能大。

## 3.3 算法流程总结

经过上面的讨论, 算法流程如图4所示。

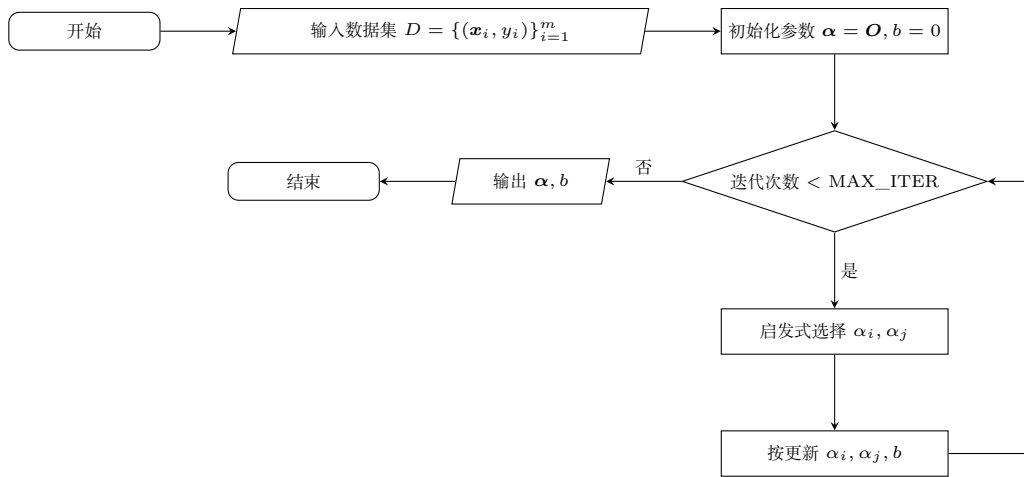


图 4: 算法流程

## 4 核心代码分析

### 4.1 数据定义

代码中定义最大迭代次数，用矩阵  $\mathbf{X}$  表示西瓜数据集 3.0 $\alpha$  的数据部分，每个样本有密度、含糖率两个维度，用向量  $y$  表示每个西瓜的真实标记。初始化  $\alpha = \mathbf{O}$ ， $b = 0$ 。定义常数  $C = 1$ ， $\varepsilon = 0.00001$ 。

```

1 max_iter = 50
2 m = 17
3 X = np.array([[0.697, 0.460], [0.774, 0.376], [0.634, 0.264], [0.608, 0.318], [0.556, 0.215], [0.403, 0.237],
    [0.481, 0.149], [0.437, 0.211], [0.666, 0.091], [0.243, 0.267], [0.245, 0.057], [0.343, 0.099], [0.639,
    0.161], [0.657, 0.198], [0.360, 0.370], [0.593, 0.042], [0.719, 0.103]])
4 # 真实标记
5 y = np.array([1, 1, 1, 1, 1, 1, 1, 1, -1, -1, -1, -1, -1, -1, -1, -1])
6 # 初始化
7 alpha = np.zeros(m)
8 b = 0
9 # 常数取值
10 epsilon = 0.00001
11 C = 1

```

### 4.2 定义核函数 $\kappa$

kappa 函数提供六种核函数，包括线性核、高斯核、多项式核、拉普拉斯核、Sigmoid 核、RBF 核，函数默认使用高斯核。

```

1 def kappa(x1, x2, func='gauss'):
2     if func == 'gauss':
3         sigma = 0.11855
4         return np.exp(-(np.linalg.norm(x1 - x2) ** 2) / (2 * sigma ** 2))
5     elif func == 'linear':
6         return np.dot(x1, x2)
7     elif func == 'polynomial':
8         d = 5
9         return np.dot(x1, x2) ** d
10    elif func == 'laplace':
11        sigma = 0.11855
12        return np.exp(-np.linalg.norm(x1 - x2) / sigma)
13    elif func == 'sigmoid':
14        beta = 0.8
15        theta = -0.2
16        return np.tanh(beta * np.dot(x1, x2) + theta)
17    elif func == 'RBF':
18        gamma = 0.2
19        return np.exp(-gamma * np.linalg.norm(x1 - x2) ** 2)
20    else:
21        raise Exception('核函数不存在')

```

### 4.3 修剪操作

修剪操作，实际上就是让超过可行域的取值  $a$  落到可行域的边界上。

```

1 def clip(a, _L, _H):
2     if a < _L:
3         return _L
4     elif a > _H:
5         return _H
6     else:
7         return a

```

### 4.4 SMO 过程

按照图4的流程，每一次迭代都按照启发式规则选择  $\alpha_i$  和  $\alpha_j$ ：首先选择不满足式(40)的  $\alpha_i$ ，若都满足约束则随机选择；接着再选择  $\alpha \arg \max_{1 \leq j \leq m, j \neq i} |E_i - E_j|$ ，更新  $\alpha_i$ 、 $\alpha_j$ 、 $b$ 。

```

1 for _ in range(max_iter):
2     index1 = index2 = -1
3     # 寻找不满足约束的alpha_i
4     for i in range(m):
5         Ei = f(X[i]) - y[i]
6         if alpha[i] < C and y[i] * Ei < -epsilon:
7             index1 = i
8             break
9         if alpha[i] > 0 and y[i] * Ei > epsilon:
10            index1 = i
11            break
12
13     # 都满足约束则随机选择
14     if index1 == -1:
15         index1 = np.random.randint(0, m)
16
17     E1 = f(X[index1]) - y[index1]
18
19     # 寻找使更新步长最大的alpha_j
20     E2 = E1
21     for i in range(m):
22         if i == index1:
23             continue
24         Ei = f(X[i]) - y[i]
25         if np.abs(Ei - E1) > np.abs(E2 - E1):
26             E2 = Ei
27             index2 = i
28
29     alpha1_old = alpha[index1]
30     alpha2_old = alpha[index2]
31     # 更新
32     alpha2_new = alpha2_old + y[index2] * (E1 - E2) / (kappa(X[index1], X[index1]) + kappa(X[index2], X[
        index2]) - 2 * kappa(X[index1], X[index2]))

```

```

33     # 确定可行域
34     if y[index1] == y[index2]:
35         L = max(0, alpha1_old + alpha2_old - C)
36         H = min(C, alpha2_old + alpha1_old)
37     else:
38         L = max(0, alpha2_old - alpha1_old)
39         H = min(C, alpha2_old - alpha1_old + C)
40
41     # 修剪
42     alpha2_new = clip(alpha2_new, L, H)
43
44     alpha1_new = alpha1_old + y[index1] * y[index2] * (alpha2_old - alpha2_new)
45     alpha[index1] = alpha1_new
46     alpha[index2] = alpha2_new
47
48     # 更新b
49     b1 = -E1 - y[index1] * kappa(X[index1], X[index1]) * (alpha1_new - alpha1_old) - y[index2] * kappa(X[
        index1], X[index2]) * (alpha2_new - alpha2_old) + b
50     b2 = -E2 - y[index1] * kappa(X[index1], X[index2]) * (alpha1_new - alpha1_old) - y[index2] * kappa(X[
        index2], X[index2]) * (alpha2_new - alpha2_old) + b
51     if 0 < alpha1_new < C:
52         b = b1
53     elif 0 < alpha2_new < C:
54         b = b2
55     elif L != H:
56         b = (b1 + b2) / 2

```

## 4.5 结果预测

利用式(25)得到  $f(x_i)$ 。

```

1 def f(_x):
2     ans = b
3     for i in range(m):
4         ans += alpha[i] * y[i] * kappa(_x, X[i])
5     return ans

```

## 5 实验结果与分析

此次实验中用到了表1中所列举的核函数，使用不同核函数，对西瓜数据集 3.0 $\alpha$  的分类效果如图5-10所示。

基于在西瓜数据集 3.0 $\alpha$  上的分类效果，我们可以初步得到如下结论：

- 线性核、多项式核、RBF 核、Sigmoid 核的分类效果类似，都没有将好瓜与坏瓜完全区分开来；
- 高斯核、拉普拉斯核的分类效果相似，将好瓜与坏瓜完全区分开来。

如表1所示，高斯核、拉普拉斯核、RBF 核的数学形式相似，因此高斯核、拉普拉斯核的分类效果相似，RBF 核分类效果差的原因在于  $\gamma$  取值的限制。在 scikit-learn 提供的 SVM 包中，取消了  $0 < \gamma < 1$  这一限制，仅提供 RBF 核和拉普拉斯核，并默认  $\gamma = \frac{1}{n_{\text{features}} \cdot X.\text{var}()}$ ，也能取得较好的分类效果。

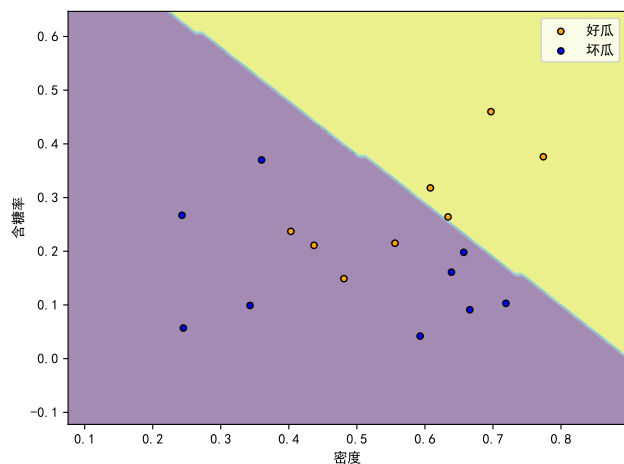


图 5: 线性核

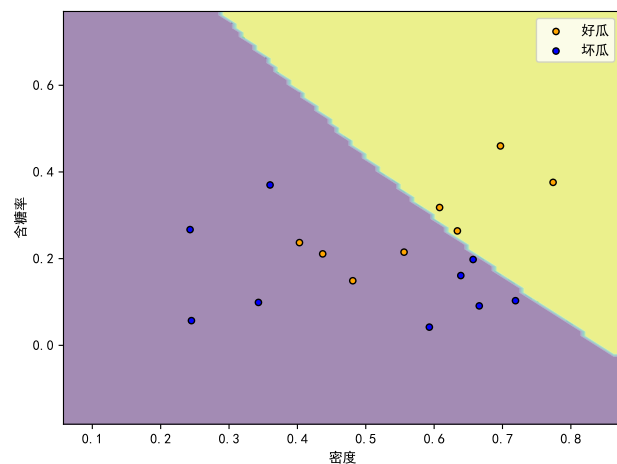


图 6: 多项式核

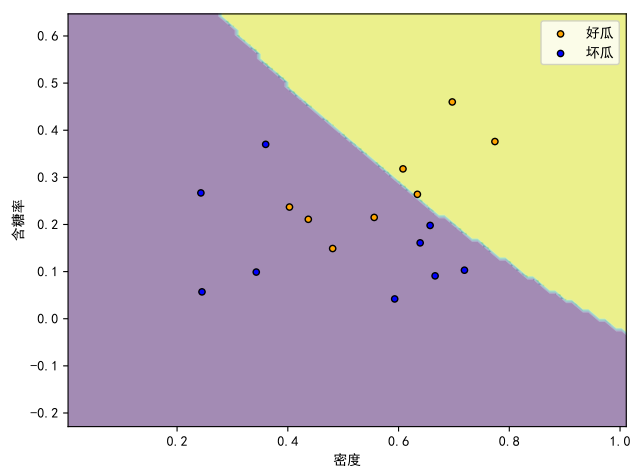


图 7: RBF 核

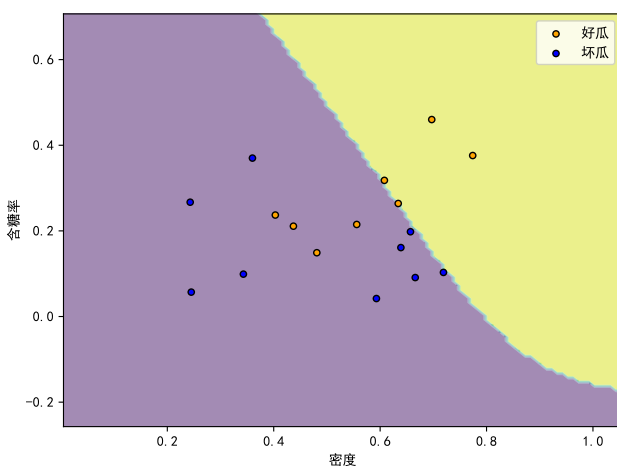


图 8: Sigmoid 核

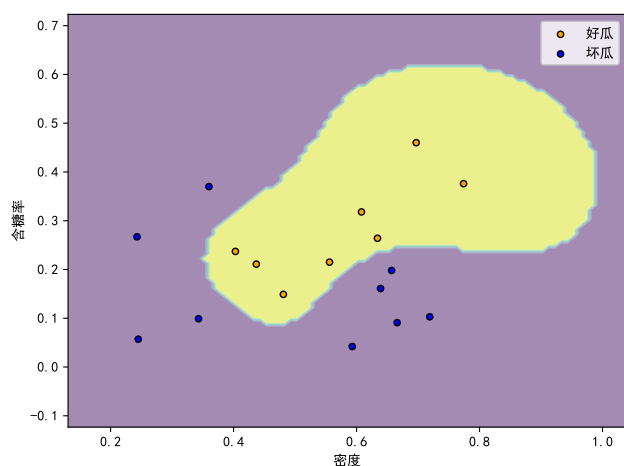


图 9: 高斯核

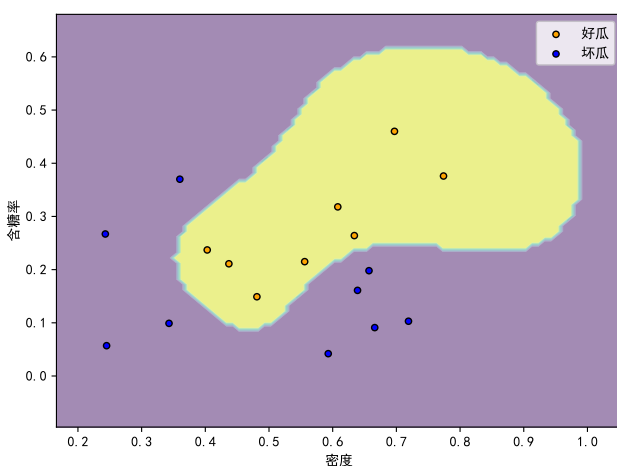


图 10: 拉普拉斯核

线性核、多项式核、Sigmoid 核表示的函数在二维空间中都是不封闭的曲线，可以看到二维空间中的分类边界就类似这些函数的曲线，仅从西瓜数据集 3.0 $\alpha$  来看，并不适合作为 SVM 的核函数。

线性核计算快、可计算得到明确的边界，但是只能解决线性可分问题。多项式核、Sigmoid 核可解决非线性问题，多项式核对于大数量级的幂数计算较慢，而 Sigmoid 核更常用于训练多层感知机神经网络。高斯

核、RBF 核、拉普拉斯核能将样本映射到无穷维空间，决策边界更为多样，分类效果好，但是难以计算清晰的分类边界。

## 6 学习收获

在这次实验过程中，我从头实现了 SVM，并对多个不同核函数的分类效果进行了比较。基于西瓜数据集 3.0 $\alpha$  上的分类效果，得出了不同核函数的优缺点和适用场景。

关于 SMO 方法、KKT 条件等知识，书上提及甚少，我因此查阅了大量资料，自己推导了 SMO 算法的诸多迭代式。实际上，在 SMO 过程中，我们还可以动态更新  $E_i$ ，而不必每一步迭代都重新计算  $E_i$ ，将来需要在代码中完善这一过程，进而提高训练 SVM 的速度。

此次实验还留下疑问，即 Sigmoid 核与多层感知机神经网络之间的联系。

## 7 参考资料

- 《机器学习》6.1,6.2,6.3,6.4；周志华；清华大学出版社
- 支持向量机 (scikit-learn 中文社区)
- sklearn.svm.SVC (scikit-learn 中文社区)
- Why does the RBF (radial basis function) kernel map into infinite dimensional space, mentioned many times in machine learning lectures?
- 【吃瓜教程】《机器学习公式详解》(南瓜书) 与西瓜书公式推导直播合集 | 第六章-支持向量机
- 【机器学习详解】SMO 算法剖析
- SEVEN MOST POPULAR SVM KERNELS
- 支持向量机 (SVM) | SMO 方法篇
- 机器学习算法实践-SVM 中的 SMO 算法
- 支持向量机原理详解 (五): KKT 条件 (Part II)
- 支持向量机原理详解 (六): 序列最小最优化 (SMO) 算法 (Part I)