

机器学习概论实验报告

PB20111699 吴骏东

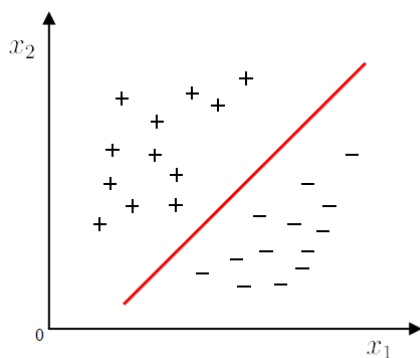
2022.10.30

在本次的实验中，我们需要自己实现支持向量机这一机器学习基本任务模型。本文将从算法原理、代码实现、性能分析等方面展开。如果对于其中部分内容有疑问或者有其他改进意见，欢迎提出 issue 或直接联系作者。

原理分析

支持向量机

支持向量机 (support vector machines, SVM) 所解决的是二分类目标中确定最大间隔超平面的问题。



如上图所示，对于二分类样本集合 $D \in (R^n \times \{-1, 1\})^m$ ，假设超平面能够将所有训练样本正确分类，即对于 $\forall (x_i, y_i) \in D$ ，有 $y_i(w^T x_i + b) > 0$ 。我们定义超平面关于样本点 (x_i, y_i) 的几何间隔为

$$\gamma_i = y_i d_i = y_i \left(\frac{w^T \cdot x_i + b}{\|w\|} \right)$$

我们希望寻找一个最优的超平面，使得所有样本点到该超平面距离“最大”。则支持向量机的优化问题可以表示为

$$\begin{aligned} & \max_{w, b} \gamma \\ & \text{s.t. } \gamma_i \geq \gamma, i = 1, 2, \dots, m \end{aligned}$$

不失一般性，令 $w = \frac{w}{\|w\| \gamma}$ ， $b = \frac{b}{\|w\| \gamma}$ ，则上面的问题等价于

$$\begin{aligned} & \max_{w, b} \gamma \\ & \text{s.t. } y_i(w x_i + b) \geq 1, i = 1, 2, \dots, m \end{aligned}$$

由于 $\max \gamma \Leftrightarrow \max \frac{1}{\|w\|} \Leftrightarrow \min \frac{1}{2} \|w\|^2$ ，所以最终我们可以将 SVM 的最大分割超平面问题转化为如下问题

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

这是一个含有不等式约束的凸二次规划问题。

软间隔 SVM

支持向量机在线性可分数据集中表现优秀，算法效率高且准确度好。但实际生活中的数据集往往是线性不可分的，并伴随有一定的噪声与空缺。SVM 对于这部分异常数据是异常敏感的：如果支持向量的选择偏离了数据集的原始特征，则 SVM 的分类结果会出现极大的误差。

因此，对于这部分问题，如果强制要求所有的样本点都满足硬间隔，可能会导致出现过拟合的问题，甚至会使决策边界发生变化。为了解决支持向量机的这一局限性，我们引入了软间隔支持向量机。

软间隔支持向量机是为了解决线性不可分问题而设计的，它允许支持向量机在一些样本上不满足约束条件（被错误分类）。为此，我们引入如下的优化目标：

$$\min_{w,b} \frac{1}{2} \|w\|^2 + \gamma \sum_{i=1}^m \ell_{0/1}(y_i(w^T x_i + b) - 1)$$

其中 γ 为正则化参数， $\ell_{0/1}$ 是 0/1 损失函数

$$\ell_{0/1}(z) \begin{cases} 1, & \text{if } z < 0; \\ 0, & \text{otherwise.} \end{cases}$$

这个想法是自然的：在最小化超平面间隔的同时，我们尽可能让样本被错误分类的情况少发生。但由于 $\ell_{0/1}$ 非凸、不连续，本实验中我们采用 Hinge 损失函数进行替代

$$\ell_{Hinge}(z) = \max(0, 1 - z)$$

于是得到了最终的优化目标

$$\min_{w,b} \frac{1}{2} \|w\|^2 + \gamma \sum_{i=1}^m \ell_{Hinge}(y_i(w^T x_i + b) - 1)$$

在本实验中，我们将采用软间隔支持向量机的两种不同求解方式，分别进行对应的实现。读者可以在后文中发现两种实现的异同之处。

代码实现

基于梯度下降的软间隔 SVM 求解

理论推导

假设数据集 X 大小满足 $X \in \mathbb{R}^{m \times n}$ ，我们引入 $\hat{w} = \begin{pmatrix} w \\ b \end{pmatrix} \in \mathbb{R}^{(n+1) \times 1}$ ， $\hat{X} = \begin{pmatrix} X & \mathbf{1} \end{pmatrix} \in \mathbb{R}^{m \times (n+1)}$ ，则 $\sum_{i=1}^m y_i(w^T x_i + b) = \sum_{i=1}^m y_i(\hat{w}^T \hat{x}_i) = y \cdot (\hat{X} \hat{w})$ 。接下来，我们引入 $\xi_i = \max(0, 1 - y_i \hat{w}^T \hat{x}_i)$ ，并记

$$\hat{y} = (\hat{y}_i), \text{ where } \hat{y}_i = \begin{cases} 0, & \text{if } \xi_i = 0 \\ y_i, & \text{if } \xi_i \neq 0 \end{cases}$$

于是目标问题可以化为

$$L = \frac{1}{2} \hat{w}^T \hat{w} + \gamma \sum_{i=1}^m \xi_i$$

$$\Rightarrow \nabla L = \hat{w} - \gamma \hat{X}^T \hat{y}$$

$$w_{new} = w_{old} - lr \times \nabla L$$

上面的矩阵梯度可能有一些复杂，读者在证明时应格外注意。

代码说明

核心部分的训练代码如下，基本上是对照上面的公式进行转述。

```
m, n = X.shape
self.w = np.zeros((n + 1, 1))

temp_1 = np.ones((m, 1))
X_hat: np.ndarray = np.c_[X, temp_1]

temp_0 = np.zeros((m, 1))
loss_list = []
y_diag = np.diag(y.reshape(-1))

for times in range(max_times):
    xi = array_max(temp_0, 1 - (y_diag @ X_hat @ self.w))
    loss = 0.5 * (self.w.T @ self.w)[0][0] + gamma * (xi.sum())

    y_bar = array_find0(xi, y)

    delta_1 = self.w - gamma * (X_hat.T @ y_bar)

    if times >= 2 and abs(loss_list[-1] - loss) < tol:
        loss_list.append(loss)
        break

    self.w = self.w - lr * delta_1
    loss_list.append(loss)
```

其中，`array_find0` 函数为向量函数，采用如下的实现

```
def find_zero(a, b):
    return 0 if a == 0 else b

def array_find0(a: np.ndarray, b: np.ndarray) -> np.ndarray:
    func_ = np.frompyfunc(find_zero, 2, 1)
    return(func_(a, b))
```

该函数可以对两个输入向量的每一维进行操作，并得到一个新的向量。

基于 SMO 的软间隔 SVM 对偶问题求解

SMO 算法的原始论文可以参考 [这里](#)。

训练一个支持向量机需要解决一个非常大的二次规划优化（QP）问题。SMO 算法将这个大的 QP 问题分解为一系列可能的最小的 QP 问题，并进行求解。算法的基本思想是：如果所有变量的解都满足最优化问题的 KKT 条件，则已经得到该最优化问题的解。否则，我们可以选择两个变量，同时固定其他变量，仅针对这两个变量构建一个 QP 问题。这样，我们通过求解两个变量的 QP 问题，能让结果不断靠近原有 QP 问题的解，并且双变量 QP 问题有着对应的解析方法。

那么应当如何选择两个变量、进行现有参数的更新呢？

SVM 对偶问题

回到我们的原始问题上：

$$\begin{aligned} \min_{w,b} \quad & \frac{1}{2} \|w\|^2 \\ \text{s.t.} \quad & y_i(w x_i + b) \geq 1, i = 1, 2, \dots, m \end{aligned}$$

采用拉格朗日乘子法，构造拉格朗日函数如下

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

其中 $\alpha_i \geq 0$ 为拉格朗日乘子。令 $\frac{\partial L}{\partial w} = 0$, $\frac{\partial L}{\partial b} = 0$ ，我们可以得到

$$\begin{aligned} w &= \sum_{i=1}^m \alpha_i y_i x_i \\ \sum_{i=1}^m \alpha_i y_i &= 0 \end{aligned}$$

代回拉格朗日函数，消去 w 和 b，我们得到

$$L(w, b, \alpha) = -\frac{1}{2} \sum_{i=1}^m \sum_{j=1}^m \alpha_i \alpha_j y_i y_j (x_i^T x_j) + \sum_{i=1}^m \alpha_i$$

由此可知，SVM 问题的对偶问题为

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

启发式变量选择法

SMO 算法在每个子问题中需要选择两个变量进行优化，并且其中至少一个变量是违反 KKT 条件的。本实验中，KKT 条件为

$$\begin{cases} \alpha_i \geq 0 \\ y_i (w^T x_i + b) - 1 \geq 0 \\ \alpha_i (y_i (w^T x_i + b) - 1) \geq 0 \end{cases}$$

我们可以先找出违反 KKT 条件最为“严重”的一系列变量，按照一定顺序存入 `index_1_list`。随后，对于 `index_1_list` 中的每个变量 α_i ，遍历所有 α_j , $j \neq i$ 使得 $|E_1 - E_2|$ 达到最大。如果不存在这样的 α_j ，则顺延至下一个 α_i 。直到 $\Delta \ell < tol$ 或所有 α_i 均满足 KKT 条件为止。

以下是寻找第一个变量的代码实现。我们优先寻找满足 $0 < \alpha_i < \gamma$ 的违反 KKT 变量，为此将其违反程度 $\times 10$ 之后存入暂存列表。这样一轮遍历之后，我们即可得到一个有序列表，对应着违反程度最为严重的一系列变量。

```
for i in range(m):
    alpha_i = self.alpha[i, :][0]
    err_i = self.err[i, :][0]
    if (0 < alpha_i < gamma and abs(err_i - 1) > epsilon):
        val = (abs(err_i) - epsilon) * 10
        # 注意这里，我们优先考虑这种情况，因此乘上了一个权重系数
        i_list.append((val, i))
    elif alpha_i == 0 and err_i < 1 - epsilon:
        val = - err_i + 1 - epsilon
        i_list.append((val, i))
    elif alpha_i >= gamma and err_i > 1 + epsilon:
        val = err_i - 1 - epsilon
        i_list.append((val, i))
```

对于第二个变量 α_j ，我们只需要根据 E_i 的正负进行判断。若 $E_i < 0$ ，则选择最大的 E_j ，否则选择最小的 E_j 即可。

```
err_dict = []
err_list = []

for i in self.err.tolist():
    err_list.append(i[0])
    for index, value in enumerate(err_list):
        err_dict.append((value, index))
    err_dict.sort(key=takefirst)

    k = 0
    if e1 > 0:
        while k < m:
            a2 = err_dict[k][1]
            if a1 != a2 and self._update_alpha(a1, a2, gamma, min_delta):
                break
            k += 1
    else:
        while k < m:
            a2 = err_dict[-1 - k][1]
            if a1 != a2 and self._update_alpha(a1, a2, gamma, min_delta):
                break
            k += 1
```

双变量二次规划问题求解

该部分内容参考了 [这篇](#) 文章。

不妨假定我们已经选定 α_1, α_2 为目标变量，其余变量保持固定。并设问题的原始可行解为 $\alpha_1^{old}, \alpha_2^{old}$ ，双变量二次规划问题的最优解为 $\alpha_1^{new}, \alpha_2^{new}$ ，且沿着约束方向未裁剪的 α_2 最优解为 $\alpha_2^{uncut-new}$ 。由于约束条件 $\sum_{i=1}^m \alpha_i y_i = 0$ 以及 $0 \leq \alpha_i \leq \gamma$ 的存在，我们应有 $low \leq \alpha_2^{uncut-new} \leq high$ 。其中

$$\begin{aligned} low &= \begin{cases} \max(0, \alpha_2^{old} - \alpha_1^{old}), & y_1 = y_2 \\ \max(0, \alpha_1^{old} + \alpha_2^{old} - \gamma), & y_1 \neq y_2 \end{cases} \\ high &= \begin{cases} \min(\gamma, \gamma + \alpha_2^{old} - \alpha_1^{old}), & y_1 = y_2 \\ \min(\gamma, \alpha_1^{old} + \alpha_2^{old}), & y_1 \neq y_2 \end{cases} \end{aligned}$$

可以证明，双变量二次规划问题沿着约束方向未经剪辑的解为

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

其中

$$E_i = \sum_{i=1}^m \alpha_i y_i x_i x_i^T + b - y_i$$

代表第 i 个样本的预测误差；而

$$\eta = x_1 x_1^T + x_2 x_2^T - 2x_1 x_2^T = \|x_1 - x_2\|^2$$

为常数。

经过剪辑后的解为

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

再代回约束式可得

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

为了提升效率所做的优化

由于 SMO 算法需要对 α 进行大量的遍历与比较，难以进行矩阵化表述，所以针对代码的效率优化是很必要的。为了减少参数计算与传递，SVM2 类中引入了如下内容作为类的成员

```
def __init__(self, X:np.ndarray, y:np.ndarray):
    m, _ = X.shape
    self.X = X
    self.y = y
    self.K = self.X @ self.X.T # 细节：提前计算了所有的 x_i^T * x_j

    self.alpha = np.zeros((m, 1))
    self.b = np.random.uniform(low=0.0, high=1.0, size=1)
    self.err = np.zeros((m, 1))
    self._update_e()
```

为了便于计算，我们引入了如下的私有方法

```

def _cut(self, low, high, a2_uncut)
# 将得到的 alpha2_new 进行裁剪

def _update_alpha(self, a1_index, a2_index, gamma, min_delta)
# 更新 alpha 的值, 若成功更新则返回 True; 如果不满足设定的条件则不会更新, 并返回 False

def _update_b(self, a1_index, a2_index, a1_new, a2_new, gamma)
# 根据 alpha_new, alpha_old 更新 b 的值

def _update_e(self):
# 更新误差 E 的值
    alpha_y = self.alpha * self.y
    self.err = self.k @ alpha_y + self.b - self.y

```

得益于提前将数据存储为类的成员, 所有的函数传参时仅需要传入变量的 index, 最小化函数的空间占用与耗时。

除此以外, 在更新 α 时, 若 $|\alpha_{old} - \alpha_{new}| < \min_delta$ 则不进行更新; 在裁剪过程中, 若 $L > H$ 则不进行更新。这样可以保证迭代的速度不至于过慢。

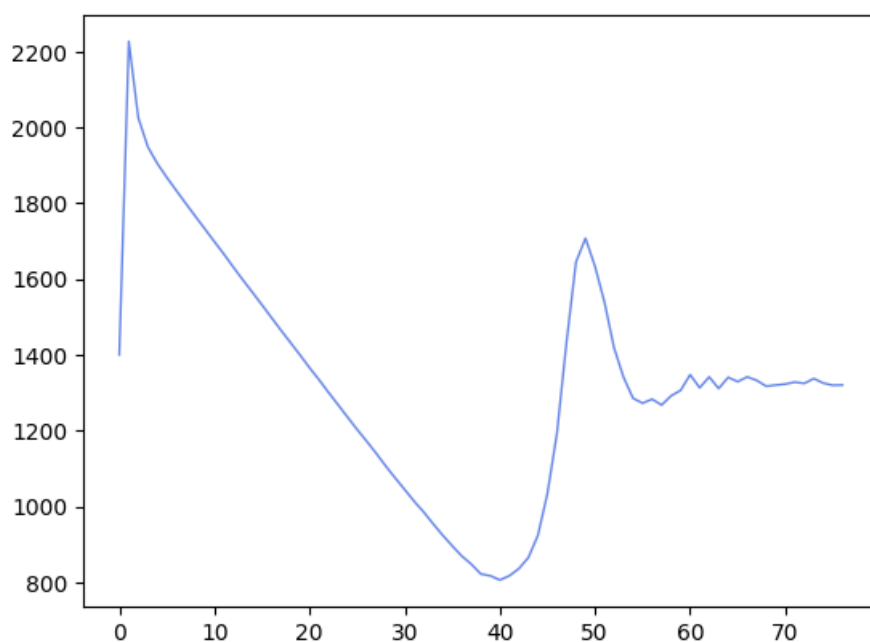
模型评价

模型 1: 基于梯度下降的软间隔 SVM 求解

由于该模型大部分的计算过程均为矩阵运算, 所以计算效率会很高。我们采用大样本数据集进行验证。

极端大样本单次验证

数据大小: 20000×50 , 错标率 0.06。训练时间: 2min35s; 模型准确率: 90.7%

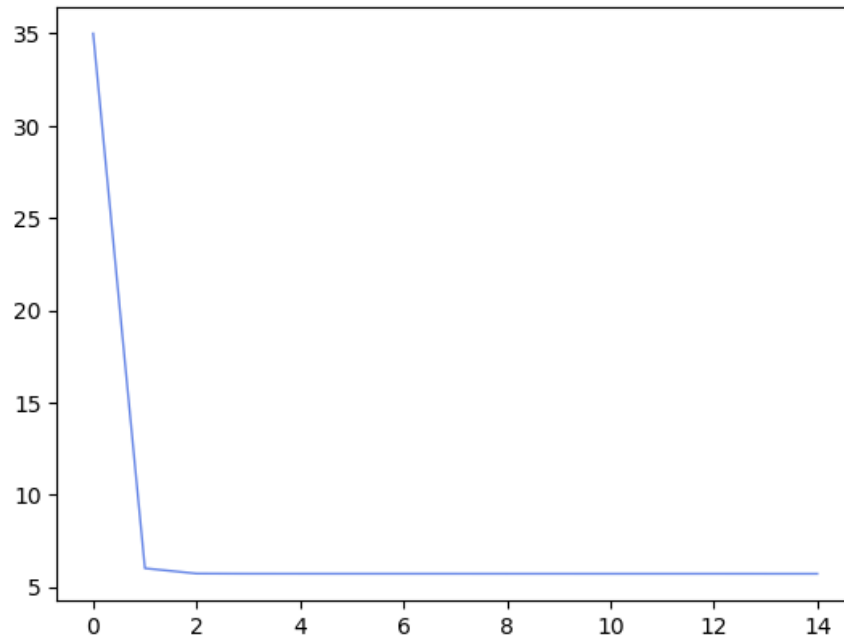


注意到损失函数曲线先是迅速下降，随后迅速上升，并逐渐在一个范围内波动。这是由于样本之中的部分误差数据引起的。为了保证训练效果，我们可以人为设置迭代次数的上界，让训练在反常上升之前即结束。

大样本单次验证

数据大小：10000 × 20，错标率 0.036。训练时间: 2.1s；模型准确率：95.1%。

参数：gamma = 0.005, lr = 0.002, tol=1e-4, max_times=100



此时的损失曲线十分正常，训练也很快就收敛了。

大样本平均验证

数据大小：10000 × 20，重复次数：100次。样本平均错标率为 0.036634，模型平均准确率为 0.9553700000000002，用时 45.7s。

```
acc, mis = model_accuracy_ave(model='1', total_time=100, dim=20, num=10000)
print("Model_{:}:Ave={}, Mis_Ave={}".format('1', acc, mis))
```

✓ 45.7s

```
Model_1:Ave=0.9553700000000002, Mis_Ave=0.036634
```

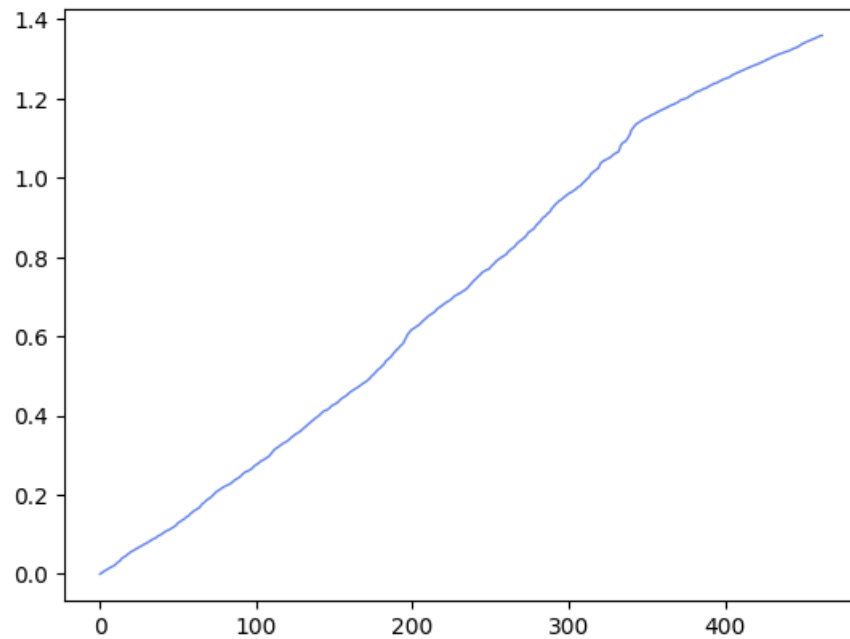
模型 2：基于 SMO 的软间隔 SVM 对偶问题求解

由于该模型的计算大部分为循环实现，且涉及启发式搜索部分，所以运算速度相对较低。我们采用小样本数据集进行验证。此外，SMO 对参数敏感，针对不同大小的数据需要更改相应的参数，否则准确率与耗时水平都会下降。

小样本单次验证

数据大小： 200×10 ，错标率 0.02。训练时间: 0.5s； 模型准确率： 95%。

参数： $\gamma = 0.1$, $\text{tol} = 1\text{e-}3$, $\text{max_times}=800$, $\text{epslion}=0.2$



损失函数值随迭代次数增加逐渐上升，符合我们的预期。然而 SMO 算法的迭代次数与训练时间都远远大于梯度下降算法。

小样本平均验证

数据大小： 200×10 ，重复次数： 50次。样本平均错标率为 0.0364，模型平均准确率为 0.91333，总用时 1min13s。

```
acc, mis = model_accuracy_ave(model='2', total_time=50, dim=10, num=200)
print("Model_{:}:Ave={}, Mis_Ave={}".format('2', acc, mis))
```

✓ 1m 13.8s

```
Model_2:Ave=0.9133333333333331, Mis_Ave=0.03640000000000001
```

模型比较

每次随机生成 50 组小样本，分别对梯度下降、SMO、sklearn.svm 三个模型进行准确率、计算时间（仅考虑模型训练与预测时间）的综合比较，结果如下

样本大小	错标率	梯度下降 准确率	梯度下降 用时	SMO 准 确率	SMO 用 时	sklearn 准确率	sklearn 用时
50×5	0.05560	0.91066	0.00146	0.90666	0.01875	0.90266	0.00062
100×10	0.04140	0.89199	0.00187	0.89199	0.04968	0.89399	0.00062
200×10	0.03930	0.91866	0.00343	0.91066	0.23750	0.91766	0.00156
500×10	0.02062	0.92720	0.02593	0.93600	6.37437	0.94013	0.02062
1000×10	0.03630	0.92919	0.06968	0.94120	49.47437	0.94660	0.05406

可以发现，梯度下降的速度与 sklearn 基本差不多，而 SMO 算法在数据规模增大的时候会急速变慢，这是大量循环与判断导致的。在准确率方面，梯度下降的准确率相对弱于 SMO 和 sklearn，且随着样本规模增大，三者准确率都有所上升。

SMO 算法的效率不尽人意。或许我们可以优化一下变量选择的策略，例如对 α_j 采用随机选择，或者限定 α_j 的选择次数等。本实验中不再进行调整。

附录

SMO 算法的一些参考参数

数据大小	γ	tol	max_times	ϵ
50×5	0.04	1e-4	500	0.45
100×10	0.04	1e-4	500	0.45
200×10	0.05	1e-4	2000	0.45
500×10	0.05	1e-4	5000	0.25
1000×10	0.05	1e-4	5000	0.25

一些辅助功能函数

```
def random_Split_data(X: np.ndarray, y: np.ndarray, rate = 0.7, random_seed: int = -1)
# 根据设定的比例进行数据集随机划分

def show(times, loss, color = '#4169E1', start=0, end=2000)
# 画图函数，范围[start, end]
# 需要 import matplotlib.pyplot as plt

def model_cmp(y_pre: np.ndarray, y_test: np.ndarray)
# 准确率比较函数

def model_accuracy_ave(model: str = '1', dim = 20, num = 10000,
                        devide_rate = 0.7, total_time = 50, ifsilent = True)
# 单模型多次训练平均效果评价函数

def model_accuracy_cmp(dim = 20, num = 10000, devide_rate = 0.7,
```

```
total_time = 50, ifsilent = True)  
# 多模型多次训练效果横向评价函数
```