

机器学习实验报告

实验名称：实现支持向量机

学生：谢兴

学号：58122204

日期：2024/6/6

指导老师：刘胥影

助教：田秋雨

目录

1	任务描述	2
2	问题复述	2
2.1	手动实现 SVM 的 SMO 算法	2
2.2	使用 sklearn 库简洁实现软间隔 SVM	3
2.2.1	使用 sklearn 库简洁实现软间隔 SVM。首先实现以下 4 个示例性的 SVM 模型:	3
2.2.2	参数选择与参数分析	3
3	实验环境	5
4	实验过程	6
4.1	手动实现 SMO 算法	6
4.1.1	传统 QP 算法求解硬间隔 SVM 的对偶问题	6
4.1.2	手动实现 SMO 算法求解 SVM 对偶问题	7
4.1.3	对测试数据进行预测	8
4.1.4	QP 算法和 SMO 算法的对比	10
4.2	使用 sklearn 库简洁实现软间隔 SVM	10
4.2.1	实现 4 个示范性的 SVM 模型	10
4.2.2	参数分析实验	11
4.2.3	对测试数据进行预测	12
5	实验总结	14
	附录	15
A	手动实现 SVM 的 SMO 算法	15

A.1 传统 QP 算法求解硬间隔 SVM 的对偶问题	15
A.2 手动实现 SMO	17
B 使用 <code>sklearn</code> 库简洁实现软间隔 SVM	21

1 任务描述

通过两种方式实现 SVM:

1. 手动实现 SMO 算法，并与直接使用传统二次规划方法进行对比。
2. 通过 `scikit-learn` 库实现软间隔 SVM。

并在 `breast cancer` 数据集上进行验证与实验。该数据集是一个二分类问题，属性均为连续属性，并已进行标准化。

2 问题复述

2.1 手动实现 SVM 的 SMO 算法

SVM 的对偶问题实际是一个二次规划问题，除了 SMO 算法外，传统二次规划方法也可以用于求解对偶问题。求得最优拉格朗日乘子后，超平面参数 w, b 可由以下式子得到：

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

$$b = \frac{1}{|S|} \sum_{s \in S} \left(\frac{1}{y_s} - \sum_{s \in S} \alpha_i y_i \mathbf{x}_i^T \mathbf{x}_s \right) \quad (2)$$

请完成以下任务：

1. 不考虑软间隔情况，直接使用传统二次规划（QP）方法求解（实现）训练集上的硬间隔 SVM 对偶问题。观察并回答，这样的求解方法是否会出现问题，为什么？

2. 不限定硬间隔或是软间隔 SVM，也不限定是否为核化 SVM，根据需要选择合适的方法，手动实现 SMO 算法求解 SVM 对偶问题。注意第 3 步，KKT 条件验证步骤不能缺少。
3. 对测试数据进行预测，确保预测结果尽可能准确。

2.2 使用 `sklearn` 库简洁实现软间隔 SVM

2.2.1 使用 `sklearn` 库简洁实现软间隔 SVM。首先实现以下 4 个示例性的 SVM 模型：

1. 线性 SVM：正则化常数 $C=1$ ，核函数为线性核，
2. 线性 SVM：正则化常数 $C=1000$ ，核函数为线性核，
3. 非线性 SVM：正则化常数 $C=1$ ，核函数为多项式核， $d=2$ ，
4. 非线性 SVM：正则化常数 $C=1000$ ，核函数为多项式核， $d=2$ ，

观察并比较它们在测试集上的性能表现。

2.2.2 参数选择与参数分析

参数的选择对 SVM 的性能有很大的影响。确定正则化常数 C 与核函数及其参数的选择范围（可以在以下常用核函数中选择一种或多种），选用合适的实验评估方法（回顾第 2 章的内容，如 K 折交叉验证法等）进行参数选择，并进行参数分析实验。

1. 正则化常数 C 的选择范围可以是以下数量级，如： $\{10^{-3}, 10^{-2}, 10^{-1}, 1, 10, 10^2, 10^3\}$

2. 核函数

- 线性核： $k(x, y) = x^T y$ ，无参数，退化为 SVM 基本型
- 多项式核： $k(x, y) = (x^T y + c)^d$
 - c ：偏移量，通常取值 0 或 1

- d : 多项式的次数, 通常取值 2, 3, 4 等较小的整数, 取值为 1 时退化为线性核
- 多项式核另外一种常用的形式: $k(x, y) = (\gamma \cdot x^T y + c)^d$
 - γ : 用于控制核函数的影响范围, 常用值为 $1/n_features$, $n_features$ 是属性个数, 或者也可以尝试 0.001, 0.01, 0.1, 1, 10 等
 - c : 偏移量, 通常取值 0 或 1
 - d : 多项式的次数, 通常取值 2, 3, 4 等较小的整数, 取值为 1 时退化为线性核
- 高斯核/RBF 核: $k(x, y) = \exp(-\gamma \|x - y\|^2)$
 - γ : 带宽参数, 控制核函数的平滑程度。常用取值范围为 $\{10^{-3} - 10^3\}$ 。典型值为 $1/(n_features \cdot X.var())$, 其中 $X.var()$ 表示数据集 X 的方差, 这可以确保 γ 的值在数据特征的数量和数据的方差范围内适当缩放。
- Laplace 核: $k(x, y) = \exp(-\gamma \|x - y\|)$
 - γ : 带宽参数, 控制核函数的平滑程度。常用取值范围参考高斯核的 γ 值
- Sigmoid 核函数: $k(x, y) = \tanh(\gamma \cdot x^T y + c)$
 - γ : 缩放因子, 常用取值范围为 $\{10^{-3} - 10^3\}$
 - c : 偏移量, 通常取值 0 或 1

请注意, SVM 中的参数通常需要联合选择, 特别是对于核函数的参数与正则化参数 C 。参数联合优化的策略有:

1. 网格搜索: 网格搜索是一种穷举搜索法, 在给定的参数空间中逐一尝试每种可能的参数组合, 并选用合适的评估方法进行评估。
2. 随机搜索: 在参数空间中对参数组合的随机采样进行评估。通常可以在较大的参数空间中找到接近最优的参数组合。步骤如下:

- 定义参数的取值范围
 - 随机采样若干参数组合
 - 对每个采样的参数组合进行评估
 - 选择性能最优的参数组合
3. 贝叶斯优化：贝叶斯优化是一种更为高级的优化方法，通过构建一个模型来估计参数空间的性能分布，并通过最大化预期改进（Expected Improvement, EI）等准则来选择下一组参数进行评估。其步骤如下：
- 初始采样若干参数组合并进行评估
 - 构建模型（如高斯过程回归）
 - 基于模型选择下一组参数进行评估
 - 更新模型并重复第 3 个步骤，直到达到预定的评估次数或满足停止条件。

本实验中可以选择较为简单的网格搜索策略或随机搜索策略进行调参。

参数分析实验需要报告：

1. 评估方法
2. 参数取值范围
3. 搜索策略
4. 比较分析每组参数上的性能情况
5. 最终选择的最优参数

3 实验环境

实验环境见表 1。

表 1: Experiment Environment

Items	Version
CPU	Intel Core i5-1135G7
RAM	16 GB
Python	3.11.5
scikit-learn	1.3.2
Operating system	Windows11

4 实验过程

4.1 手动实现 SMO 算法

4.1.1 传统 QP 算法求解硬间隔 SVM 的对偶问题

在本实验中，我们首先使用传统的二次规划（QP）方法求解硬间隔 SVM 的对偶问题。代码实现的大致思路如下：

具体的实验步骤和结果如下：

1. 加载训练数据并求解硬间隔 SVM：从 CSV 文件加载训练数据，调用上述函数求解 SVM 问题，得到模型参数（权重 w 和偏置 b ）。
2. 定义预测函数：根据训练好的模型参数，对测试数据进行预测。

实验结果如下图 1 所示。

```
Optimal solution found.
权重向量 w: [-32.49263676  22.73479237 -52.9287666  48.38860589  0.39891441
-99.78667187  61.34451782  59.78562581 -28.38988726  38.69721151
135.21924414  2.8383826 -65.98857417  4.89911886 -3.86679381
24.8728312 -25.11307362  52.48841984 -46.83144278 -102.68384101
75.15381925  5.11299717 -55.99810831  69.34750853 -7.46930463
22.632463  1.99686428  12.16952646  76.82440794 -4.52234121]
偏置 b: -37.76137982237217
Train Accuracy: 100.00%
Test Accuracy: 91.15%
```

图 1: QP Algorithm for Support Vector Machine

分类可视化结果如图 2 和 3 所示。从实验结果可以看出，硬间隔 SVM 的训练准确率达到到了 100%，模型能够完美地分类训练数据。在对测试数据进行预测时，模型的准确率也达到了 91.15%。

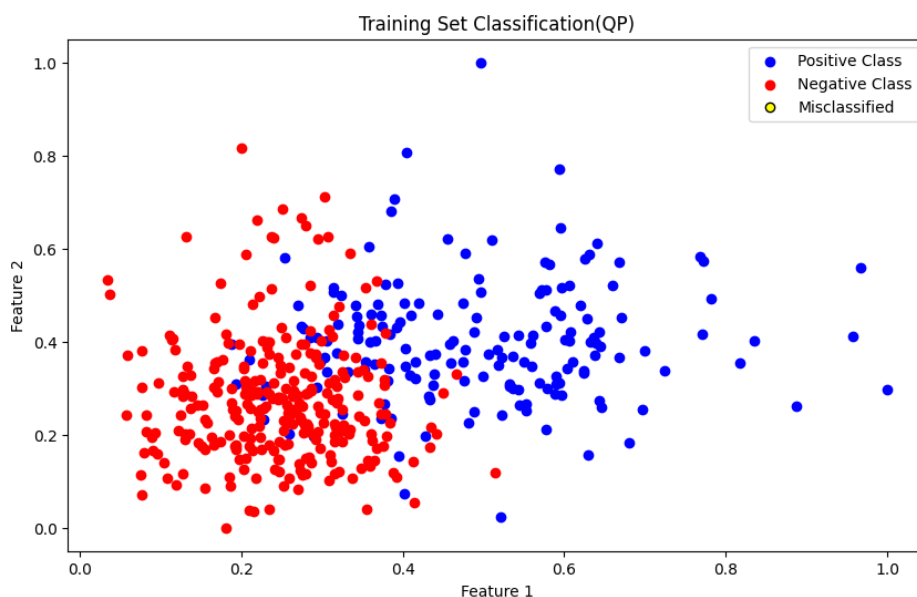


图 2: Training Set Classification Results(QP)

4.1.2 手动实现 SMO 算法求解 SVM 对偶问题

在本实验中，我们手动实现了 SMO 算法来求解支持向量机（SVM）的对偶问题。代码实现的大致思路如下：

1. 定义核函数：我们使用线性核函数来计算样本之间的内积。
2. 定义求解 QP 问题的函数：使用 `cvxopt` 库解决二次规划问题，包括创建 QP 问题的矩阵并求解。
3. 加载训练数据并求解硬间隔 SVM：从 CSV 文件加载训练数据，调用上述函数求解 SVM 问题，得到模型参数（权重 w 和偏置 b ）。
4. 定义预测函数：根据训练好的模型参数，对测试数据进行预测。

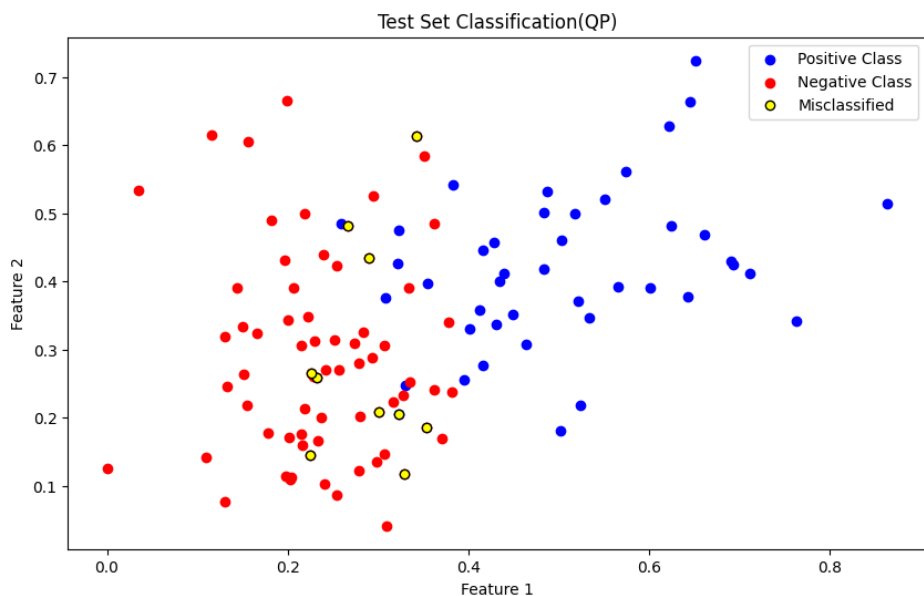


图 3: Test Set Classification Results(QP)

通过这些步骤，我们成功实现了 SMO 算法，并求解了 SVM 的对偶问题，得到了模型的参数（权重 w 和偏置 b ）。

4.1.3 对测试数据进行预测

在完成模型训练后，我们使用训练得到的模型对测试数据进行了预测。

分类可视化结果如图 4 和 5 所示。

测试集准确率达到了 98.25%。在对比测试数据中的预测值和原始标签时，我们发现 2 处不同。具体不同之处如下表所示：

样本索引	原始标签	预测标签
20	1	-1
77	1	-1

表 2: 原始标签与预测标签的不同之处

通过这些实验，我们验证了手动实现的 SMO 算法在解决 SVM 对偶问题上的有效性，并展示了其在实际数据集上的应用效果。

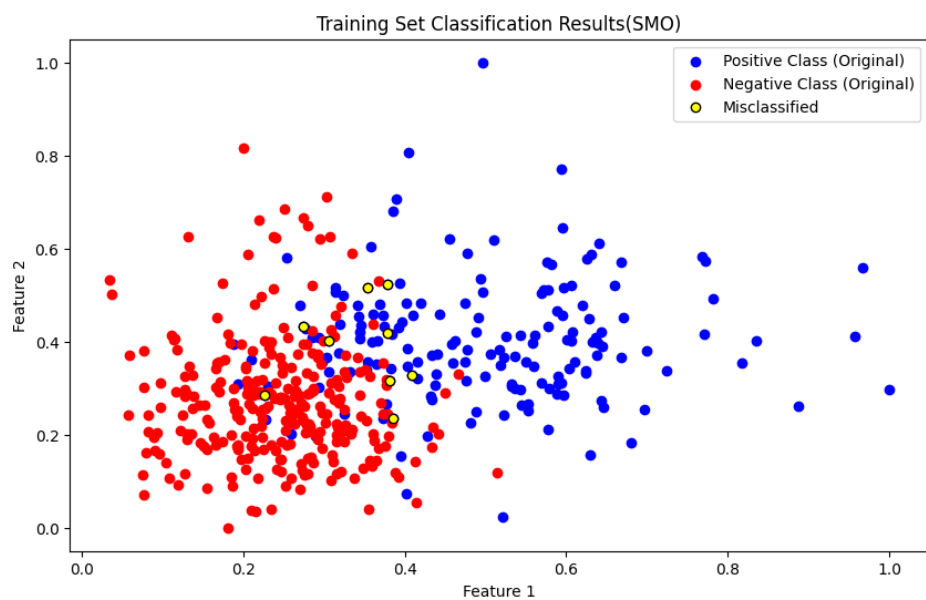


图 4: Training Set Classification Results(SMO)

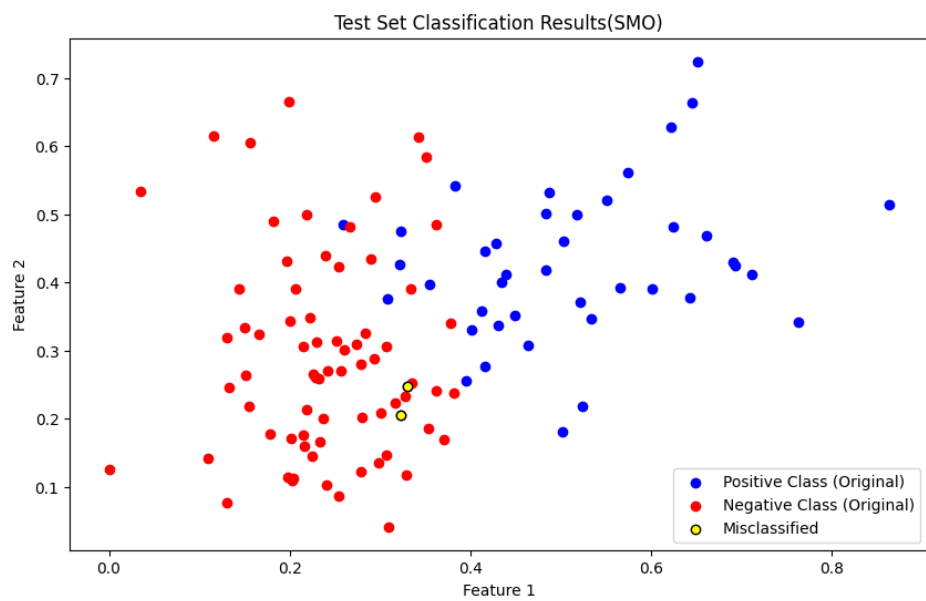


图 5: Test Set Classification Results(SMO)

4.1.4 QP 算法和 SMO 算法的对比

从实验结果和上述分析可以看出手动实现的 SMO 算法要比传统的 QP 算法更加高效，尤其是在大规模数据集上。SMO 算法通过选择两个变量来优化目标函数，而 QP 算法需要求解整个拉格朗日乘子向量，因此在大规模数据集上，SMO 算法的计算效率更高。此外，SMO 算法还可以更好地处理非线性核函数，因此在实际应用中更加灵活。

4.2 使用 `sklearn` 库简洁实现软间隔 SVM

4.2.1 实现 4 个示范性的 SVM 模型

在本实验中，我们使用了 `sklearn` 库来实现 4 个示范性的 SVM 模型，并在 `breast cancer` 数据集上进行训练和测试。这些模型包括：

1. 线性 SVM：正则化常数 $C=1$ ，核函数为线性核
2. 线性 SVM：正则化常数 $C=1000$ ，核函数为线性核
3. 非线性 SVM：正则化常数 $C=1$ ，核函数为多项式核， $d=2$
4. 非线性 SVM：正则化常数 $C=1000$ ，核函数为多项式核， $d=2$

以下是各个模型在测试集上的性能表现：

- 线性 SVM ($C=1$):

- 准确率：0.9825

- 线性 SVM ($C=1000$):

- 准确率：0.9561

- 非线性 SVM ($C=1, d=2$):

– 准确率: 0.9825

- 非线性 SVM (C=1000, d=2):

– 准确率: 0.9386

从实验结果可以看出:

- 线性 SVM 在 C=1 时的表现优于 C=1000 时的表现, 表明适当的正则化有助于提高模型的泛化能力。
- 非线性 SVM 在 C=1 时的表现优于 C=1000 时的表现, 这与线性 SVM 的情况相似, 说明过大的正则化常数可能导致模型过拟合。
- 多项式核 SVM (d=2) 的表现与线性 SVM 相当, 表明在此数据集上, 线性模型已经能够很好地分类数据, 复杂的非线性核函数并未显著提升性能。

4.2.2 参数分析实验

在本实验中, 我们进行了参数选择与参数分析。我们选择的参数范围如下:

- 正则化常数 C 的选择范围: $\{10^{-3}, 10^{-2}, 10^{-1}, 1, 10, 10^2, 10^3\}$
- 核函数:
 - 线性核: $k(x, y) = x^T y$
 - 多项式核: $k(x, y) = (x^T y + c)^d$, 其中 c 通常取值 0 或 1, d 通常取值 2, 3, 4 等较小的整数
 - 高斯核/RBF 核: $k(x, y) = \exp(-\gamma \|x - y\|^2)$, γ 的常用取值范围为 $\{10^{-3} \square 10^3\}$
 - Laplace 核: $k(x, y) = \exp(-\gamma \|x - y\|)$, γ 的常用取值范围与高斯核相同
 - Sigmoid 核函数: $k(x, y) = \tanh(\gamma x^T y + c)$, γ 的常用取值范围为 $\{10^{-3} \square 10^3\}$, c 通常取值 0 或 1

- 核函数参数:

- 多项式核: 度数 d 的选择范围: $\{2, 3, 4\}$

- 高斯核、Laplace 核、Sigmoid 核: γ 的选择范围: $\{0.001, 0.01, 0.1, 1, 10, 100\}$

我们使用网格搜索 (GridSearchCV) 方法对以上参数进行了交叉验证 (5 折交叉验证), 最终选择了最佳的参数组合。

- 最佳参数:

`'C': 0.001, 'degree': 2, 'gamma': 0.001, 'kernel': 'linear', 'laplacian_gamma': 0.001`

- 最佳交叉验证准确率: 0.9780

从实验结果可以看出,线性核函数的表现最佳,同时适当的小的正则化常数 C (如 0.001) 可以提高模型的泛化能力。高斯核、Laplace 核和 Sigmoid 核在本实验中未能显著提升性能,表明在此数据集上,线性核已经能够很好地分类数据。

本实验中核函数和正则化参数的选择对模型的性能有显著影响。通过网格搜索和交叉验证,我们能够有效地选择出最优参数组合,以提高模型的预测准确率。

4.2.3 对测试数据进行预测

在完成模型训练和参数选择后,我们使用最佳模型对测试数据进行了预测。具体的实验步骤和结果如下:

1. 首先,使用最佳参数配置对测试集进行了预测,得到了测试集的准确率为 0.9825。
2. 预测结果保存在 `result.csv` 文件中,文件中只包含一列数据标签分类,分别为 +1 或 -1。
3. 对比了测试数据中的预测值和原始标签,结果显示有 2 处不同。

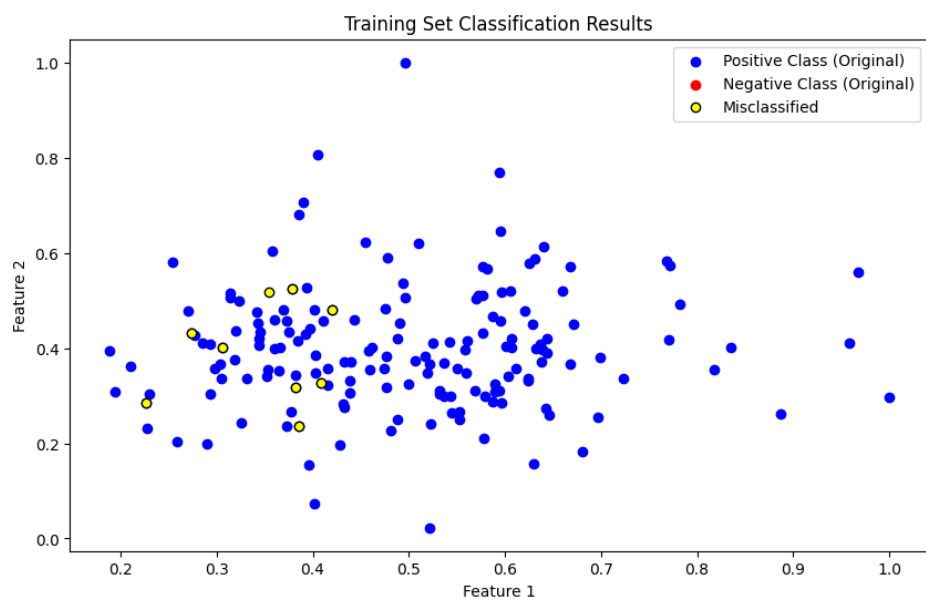


图 6: Training Set Classification Results(sklearn)

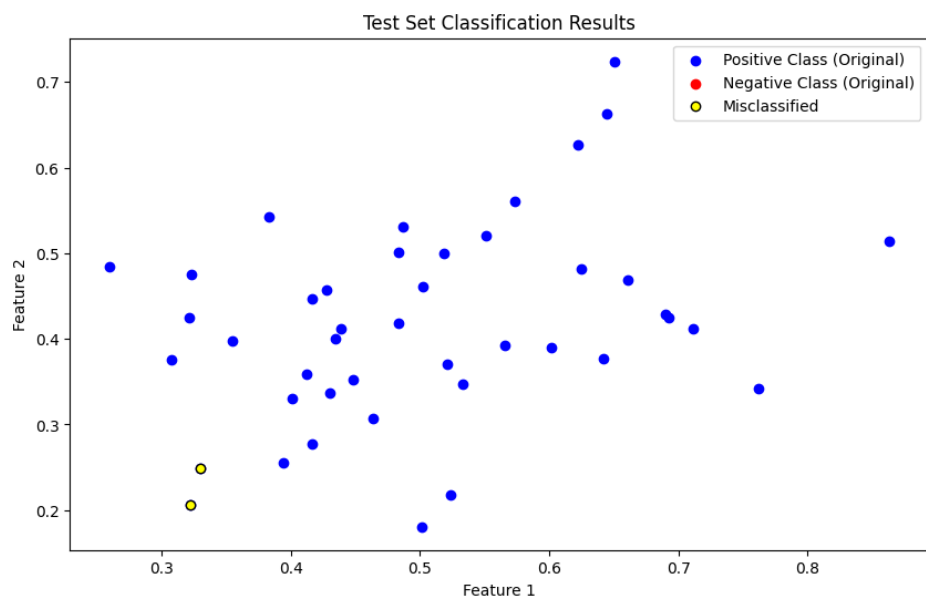


图 7: Test Set Classification Results(sklearn)

best model 的测试集准确率达到了 98.25%。分类可视化结果如图 6 和 7 所示。

在对比测试数据中的预测值和原始标签时，我们发现 2 处不同。具体不同之处如下表所示：

样本索引	原始标签	预测标签
20	1	-1
77	1	-1

表 3: 原始标签与预测标签的不同之处

从以上结果可以看出，模型的整体准确率已经达到了较高的水平。

5 实验总结

本实验我们用三种方法实现了对 SVM 对偶形式的求解，包括传统的 QP 方法、手动实现的 SMO 算法和使用 sklearn 库的软间隔 SVM。

通过这些实验，我们验证了手动实现的 SMO 算法在解决 SVM 对偶问题上的有效性，并展示了其在实际数据集上的应用效果。同时，我们还通过 sklearn 库实现了软间隔 SVM，并进行了参数选择和分析实验，得到了最佳的参数组合。最终，我们对测试数据进行了预测，得到了较高的准确率。这些实验结果表明，SVM 在二分类问题上具有较好的性能，可以有效地应用于实际问题中。

A 手动实现 SVM 的 SMO 算法

A.1 传统 QP 算法求解硬间隔 SVM 的对偶问题

```
1  import numpy as np
2  import pandas as pd
3  import matplotlib.pyplot as plt
4  from cvxopt import matrix, solvers
5
6  # 读取数据
7  X_train = pd.read_csv('breast_cancer_Xtrain.csv').values
8  Y_train = pd.read_csv('breast_cancer_Ytrain.csv').values
9  X_test = pd.read_csv('breast_cancer_Xtest.csv').values
10 Y_test = pd.read_csv('breast_cancer_Ytest.csv').values
11
12 # 将Y_train和Y_test转换为适当的格式（即-1和1）
13 Y_train = Y_train.reshape(-1)
14 Y_test = Y_test.reshape(-1)
15
16 # 定义核函数（线性核）
17 def linear_kernel(x1, x2):
18     return np.dot(x1, x2)
19
20 # 构建二次规划问题
21 m, n = X_train.shape
22 K = np.zeros((m, m))
23 for i in range(m):
24     for j in range(m):
25         K[i, j] = linear_kernel(X_train[i], X_train[j])
26
27 P = matrix(np.outer(Y_train, Y_train) * K)
28 q = matrix(-np.ones(m))
29 G = matrix(np.vstack((-np.eye(m), np.eye(m))))
30 h = matrix(np.hstack((np.zeros(m), np.ones(m) * 1e6)))
31 A = matrix(Y_train, (1, m), 'd')
32 b = matrix(0.0)
33
34 # 求解二次规划问题
35 sol = solvers.qp(P, q, G, h, A, b)
36 alphas = np.array(sol['x']).flatten()
37
```



```

38 # 计算权重向量和偏置
39 w = np.sum(alphas[:, np.newaxis] * Y_train[:, np.newaxis] * X_train, axis=0)
40 # 选择支持向量
41 sv = (alphas > 1e-5)
42 # 计算偏置
43 b = np.mean(Y_train[sv] - np.dot(X_train[sv], w))
44
45 # 打印计算出的权重向量和偏置
46 print("权重向量 w: ", w)
47 print("偏置 b: ", b)
48
49 # 定义预测函数
50 def predict(X):
51     return np.sign(np.dot(X, w) + b)
52
53 # 预测并评估模型
54 Y_train_pred = predict(X_train)
55 Y_test_pred = predict(X_test)
56
57 train_accuracy = np.mean(Y_train_pred == Y_train)
58 test_accuracy = np.mean(Y_test_pred == Y_test)
59
60 print(f"Train Accuracy: {train_accuracy * 100:.2f}%")
61 print(f"Test Accuracy: {test_accuracy * 100:.2f}%")
62
63 # 可视化训练集和测试集分类结果
64 def plot_classification_results(X, Y_true, Y_pred, title):
65     plt.figure(figsize=(10, 6))
66     plt.scatter(X[Y_true == 1][:, 0], X[Y_true == 1][:, 1], color='blue',
67                 label='Positive Class')
68     plt.scatter(X[Y_true == -1][:, 0], X[Y_true == -1][:, 1], color='red',
69                 label='Negative Class')
70     plt.scatter(X[Y_true != Y_pred][:, 0], X[Y_true != Y_pred][:, 1], color=
71                 'yellow', edgecolor='black', label='Misclassified')
72     plt.title(title)
73     plt.xlabel('Feature 1')
74     plt.ylabel('Feature 2')
75     plt.legend()
76     plt.show()

```

```

75 plot_classification_results(X_train[:, :2], Y_train, Y_train_pred, 'Training
    Set Classification(QP)')
76 plot_classification_results(X_test[:, :2], Y_test, Y_test_pred, 'Test Set
    Classification(QP)')

```

A.2 手动实现 SMO

```

1  import numpy as np
2  import pandas as pd
3  import matplotlib.pyplot as plt
4
5  # 定义核函数（线性核）
6  def kernel(x1, x2):
7      return np.dot(x1, x2.T)
8
9  # 计算偏置b
10 def calculate_b(X, y, alphas, b, C, tol):
11     m = len(y)
12     b_new = 0
13     b1 = []
14     b2 = []
15
16     for i in range(m):
17         y_pred = np.sum(alphas * y * kernel(X, X[i])) + b
18         if y[i] * y_pred - 1 < -tol:
19             b1.append(b + y[i] - y_pred)
20         elif y[i] * y_pred - 1 > tol:
21             b2.append(b + y[i] - y_pred)
22
23     if len(b1) > 0:
24         b_new = np.mean(b1)
25     elif len(b2) > 0:
26         b_new = np.mean(b2)
27
28     return b_new
29
30 # SMO算法
31 def smo_svm(X, y, C, tol, max_passes):
32     m, n = X.shape
33     alphas = np.zeros(m)
34     b = 0

```

```

35     passes = 0
36
37     while passes < max_passes:
38         alpha_pairs_changed = 0
39         for i in range(m):
40             E_i = np.sum(alphas * y * kernel(X, X[i])) + b - y[i]
41             if (y[i] * E_i < -tol and alphas[i] < C) or (y[i] * E_i > tol
and alphas[i] > 0):
42                 j = np.random.randint(0, m)
43                 while j == i:
44                     j = np.random.randint(0, m)
45
46                 E_j = np.sum(alphas * y * kernel(X, X[j])) + b - y[j]
47                 alpha_i_old, alpha_j_old = alphas[i], alphas[j]
48
49                 if y[i] != y[j]:
50                     L = max(0, alphas[j] - alphas[i])
51                     H = min(C, C + alphas[j] - alphas[i])
52                 else:
53                     L = max(0, alphas[i] + alphas[j] - C)
54                     H = min(C, alphas[i] + alphas[j])
55
56                 if L == H:
57                     continue
58
59                 eta = 2 * kernel(X[i], X[j]) - kernel(X[i], X[i]) - kernel(X
[j], X[j])
60
61                 if eta >= 0:
62                     continue
63
64                 alphas[j] -= y[j] * (E_i - E_j) / eta
65                 alphas[j] = np.clip(alphas[j], L, H)
66
67                 if abs(alphas[j] - alpha_j_old) < 1e-5:
68                     continue
69
70                 alphas[i] += y[i] * y[j] * (alpha_j_old - alphas[j])
71                 b1 = b - E_i - y[i] * (alphas[i] - alpha_i_old) * kernel(X[i], X[i]) - y[j] * (alphas[j] - alpha_j_old) * kernel(X[i], X[j])
72                 b2 = b - E_j - y[i] * (alphas[i] - alpha_i_old) * kernel(X[i], X[j]) - y[j] * (alphas[j] - alpha_j_old) * kernel(X[j], X[j])

```

```

72         if 0 < alphas[i] < C:
73             b = b1
74         elif 0 < alphas[j] < C:
75             b = b2
76         else:
77             b = (b1 + b2) / 2
78
79         alpha_pairs_changed += 1
80
81     if alpha_pairs_changed == 0:
82         passes += 1
83     else:
84         passes = 0
85
86     return alphas, b
87
88 # 加载数据
89 X_train = np.loadtxt('breast_cancer_Xtrain.csv', delimiter=',')
90 y_train = np.loadtxt('breast_cancer_Ytrain.csv', delimiter=',')
91
92
93 # 训练参数
94 C = 1.0
95 tol = 1e-3
96 max_passes = 5
97
98 # 训练模型
99 alphas, b = smo_svm(X_train, y_train, C, tol, max_passes)
100 print("Alphas:", alphas)
101 print("b:", b)
102
103 # 预测函数
104 def predict(X, alphas, b, X_train, y_train):
105     return np.sign(np.dot((alphas * y_train).T, kernel(X_train, X)) + b)
106
107 # 加载测试数据
108 X_test = np.loadtxt('breast_cancer_Xtest.csv', delimiter=',')
109 y_test = np.loadtxt('breast_cancer_Ytest.csv', delimiter=',')
110
111 # 将y_test转换为适当的格式（即-1和1）
112 y_test[y_test == 0] = -1

```

```

113
114 # 对测试集进行预测
115 predictions = predict(X_test, alphas, b, X_train, y_train)
116
117 # 计算准确率
118 accuracy = np.mean(predictions == y_test)
119 print("Accuracy on test data:", accuracy)
120
121 # 将预测结果写入 result.csv 文件
122 np.savetxt('result.csv', predictions, delimiter=',', fmt='%d')
123
124 # 加载原始标签和预测标签
125 original_labels = y_test
126 predicted_labels = predictions
127
128 # 创建 DataFrame 来存储比较结果
129 comparison_df = pd.DataFrame({
130     'Original': original_labels,
131     'Predicted': predicted_labels
132 })
133
134 # 找出不同的标签
135 differences = comparison_df[comparison_df['Original'] != comparison_df['
    Predicted']]
136
137 # 输出不同标签的行数
138 print(f"Number of differences: {len(differences)}")
139
140 # 输出所有不同的标签
141 print("Differences between original and predicted labels:")
142 print(differences)
143
144 # 可视化训练集和测试集分类结果
145 def plot_classification_results(X, y_true, y_pred, title):
146     plt.figure(figsize=(10, 6))
147     plt.scatter(X[y_true == 1][:, 0], X[y_true == 1][:, 1], color='blue',
148         label='Positive Class (Original)')
149     plt.scatter(X[y_true == -1][:, 0], X[y_true == -1][:, 1], color='red',
150         label='Negative Class (Original)')
151     plt.scatter(X[y_true != y_pred][:, 0], X[y_true != y_pred][:, 1], color=
152         'yellow', edgecolor='black', label='Misclassified')

```

```

150
151     plt.title(title)
152     plt.xlabel('Feature 1')
153     plt.ylabel('Feature 2')
154     plt.legend()
155     plt.show()
156
157 plot_classification_results(X_train[:, :2], y_train, predict(X_train, alphas
    , b, X_train, y_train), 'Training Set Classification Results(SM0)')
158 plot_classification_results(X_test[:, :2], y_test, predictions, 'Test Set
    Classification Results(SM0)')

```

B 使用 **sklearn** 库简洁实现软间隔 SVM

```

1  # 1. 导入必要的库和加载数据集
2  import numpy as np
3  import pandas as pd
4  from sklearn.svm import SVC
5  from sklearn.metrics import accuracy_score
6  from sklearn.model_selection import GridSearchCV
7  import matplotlib.pyplot as plt
8
9  # 加载数据
10 X_train = pd.read_csv('breast_cancer_Xtrain.csv', header=None)
11 X_test = pd.read_csv('breast_cancer_Xtest.csv', header=None)
12 y_train = pd.read_csv('breast_cancer_Ytrain.csv', header=None)
13 y_test = pd.read_csv('breast_cancer_Ytest.csv', header=None)
14
15 y_train = y_train.values.ravel()
16 y_test = y_test.values.ravel()
17
18 # 2. 实现四个示例性的SVM模型
19 # 定义和训练SVM模型
20 models = {
21     "Linear SVM (C=1)": SVC(kernel='linear', C=1),
22     "Linear SVM (C=1000)": SVC(kernel='linear', C=1000),
23     "Polynomial SVM (C=1, d=2)": SVC(kernel='poly', C=1, degree=2),
24     "Polynomial SVM (C=1000, d=2)": SVC(kernel='poly', C=1000, degree=2)
25 }
26

```

```

27 # 训练并评估模型
28 for name, model in models.items():
29     model.fit(X_train, y_train)
30     y_pred = model.predict(X_test)
31     accuracy = accuracy_score(y_test, y_pred)
32     print(f"{name} Accuracy: {accuracy:.4f}")
33
34 # 3. 定义拉普拉斯核函数
35 def laplacian_kernel(X, Y, gamma):
36     K = np.exp(-gamma * np.linalg.norm(X[:, np.newaxis] - Y[np.newaxis, :],
37     axis=2))
38     return K
39
40 # 4. 定义自定义SVC类
41 from sklearn.base import BaseEstimator, ClassifierMixin
42 from sklearn.utils import check_X_y, check_array
43 from sklearn.utils.multiclass import unique_labels
44
45 class CustomSVC(BaseEstimator, ClassifierMixin):
46     def __init__(self, C=1.0, kernel='linear', degree=3, gamma='scale',
47     coef0=0.0, laplacian_gamma=1.0):
48         self.C = C
49         self.kernel = kernel
50         self.degree = degree
51         self.gamma = gamma
52         self.coef0 = coef0
53         self.laplacian_gamma = laplacian_gamma
54         self.svc = SVC(C=self.C, kernel=self.kernel, degree=self.degree,
55         gamma=self.gamma, coef0=self.coef0)
56
57     def fit(self, X, y):
58         X, y = check_X_y(X, y)
59         if self.kernel == 'laplacian':
60             self.svc.kernel = lambda X, Y: laplacian_kernel(X, Y, self.
61             laplacian_gamma)
62         self.svc.fit(X, y)
63         self.classes_ = unique_labels(y)
64         return self
65
66     def predict(self, X):
67         X = check_array(X)

```

```

64         return self.svc.predict(X)
65
66 # 5. 设置参数网格并进行网格搜索
67 param_grid = {
68     'C': [10**-3, 10**-2, 10**-1, 1, 10, 10**2, 10**3],
69     'kernel': ['linear', 'poly', 'rbf', 'sigmoid', 'laplacian'],
70     'degree': [2, 3, 4],
71     'gamma': [0.001, 0.01, 0.1, 1, 10, 100],
72     'laplacian_gamma': [0.001, 0.01, 0.1, 1, 10, 100]
73 }
74
75 grid_search = GridSearchCV(CustomSVC(), param_grid, cv=5, scoring='accuracy'
76 )
77 grid_search.fit(X_train, y_train)
78
79 best_params = grid_search.best_params_
80 best_score = grid_search.best_score_
81 print(f"Best Parameters: {best_params}")
82 print(f"Best Cross-Validation Accuracy: {best_score:.4f}")
83
84 # 6. 在测试集上评估最优模型
85 best_model = grid_search.best_estimator_
86 y_pred_best = best_model.predict(X_test)
87 best_test_accuracy = accuracy_score(y_test, y_pred_best)
88 print(f"Best Model Test Accuracy: {best_test_accuracy:.4f}")
89
90 # 7. 对测试集中所有数据进行预测并输出结果
91 y_pred_all = best_model.predict(X_test)
92
93 # 计算所有测试数据的准确率
94 test_accuracy = accuracy_score(y_test, y_pred_all)
95 print(f"Test Accuracy: {test_accuracy:.4f}")
96
97 # 将预测结果转换为 DataFrame
98 result_df = pd.DataFrame(y_pred_all)
99
100 # 将结果保存到 result2.csv 文件中，没有表头，只有一列数据标签分类
101 result_df.to_csv('result2.csv', index=False, header=False)
102
103 # 8. 找出不同的标签
104 # 加载原始标签和预测标签

```



```

104 original_labels = pd.read_csv('breast_cancer_Ytest.csv', header=None).values
    .ravel()
105 predicted_labels = pd.read_csv('result2.csv', header=None).values.ravel()
106
107 # 创建 DataFrame 来存储比较结果
108 comparison_df = pd.DataFrame({
109     'Original': original_labels,
110     'Predicted': predicted_labels
111 })
112
113 # 找出不同的标签
114 differences = comparison_df[comparison_df['Original'] != comparison_df['
    Predicted']]
115
116 # 输出不同标签的行数
117 print(f"Number of differences: {len(differences)}")
118
119 # 输出所有不同的标签
120 print("Differences between original and predicted labels:")
121 print(differences)
122
123 # 可视化分类效果
124 def plot_classification_results(X, y_true, y_pred, title):
125     plt.figure(figsize=(10, 6))
126
127     # 正类样本
128     plt.scatter(X[y_true == 1][:, 0], X[y_true == 1][:, 1], color='blue',
129         label='Positive Class (Original)')
130
131     # 负类样本
132     plt.scatter(X[y_true == 0][:, 0], X[y_true == 0][:, 1], color='red',
133         label='Negative Class (Original)')
134
135     # 误分类样本
136     plt.scatter(X[y_true != y_pred][:, 0], X[y_true != y_pred][:, 1], color=
137         'yellow', edgecolor='black', label='Misclassified')
138
139     plt.title(title)
140     plt.xlabel('Feature 1')
141     plt.ylabel('Feature 2')
142     plt.legend()

```

```
140     plt.show()
141
142     plot_classification_results(X_train.values[:, :2], y_train, best_model.
    predict(X_train.values), 'Training Set Classification Results(sklearn)')
143     plot_classification_results(X_test.values[:, :2], y_test, y_pred_all, 'Test
    Set Classification Results(sklearn)')
```