# 第三章

## 一、推导软-SVM 主问题的对偶问题

一·推导 软-SVM 主问题的对偶问题

soft-SVM 主问题：$\min\ w^Tw/2 + c\sum\limits_{i=1}^{n}\epsilon_i$

约束：$y_i(w^Tx_i+b) \geq 1-\epsilon_i,\ \epsilon_i \geq 0$

利用拉格朗日乘子得到无约束优化问题

$$L(w,b,\epsilon,\alpha,\mu) = \frac{1}{2}w\cdot w + c\sum_{i=1}^{n}\epsilon_i - \sum_{i=1}^{n}\alpha_i(y_i(w\cdot x_i+b)-1+\epsilon_i) - \sum_{i=1}^{n}\mu_i\epsilon_i$$

主问题重写为 $\min\limits_{w,b,\epsilon}\ \max\limits_{\alpha,\mu}\ L(w,b,\epsilon,\alpha,\mu)$

上述主问题的对偶问题 为 $\max\limits_{\alpha,\mu}\ \min\limits_{w,b,\epsilon}\ L(w,b,\epsilon,\alpha,\mu)$

最小化上述拉格朗日：对 $w,b,\epsilon$ 求偏导，并令偏导为 $0$

$\frac{\partial L}{\partial w}=0 \Rightarrow w=\sum\limits_{i=1}^{n}\alpha_iy_ix_i$, $\frac{\partial L}{\partial b}=0 \Rightarrow 0=\sum\limits_{i=1}^{n}\alpha_iy_i$, $\frac{\partial L}{\partial \epsilon}=0 \Rightarrow c-\alpha_i-\mu_i=0$

代入 $L$ 中得 $\min\limits_{w,b,\epsilon}L(w,b,\epsilon,\alpha,\mu) = -\frac{1}{2}\sum\limits_{i=1}^{n}\sum\limits_{j=1}^{n}\alpha_i\alpha_jy_iy_jx_ix_j + \sum\limits_{i=1}^{n}\alpha_i$

代入上述对偶问题得

$$\max\limits_{\alpha}\left[-\frac{1}{2}\sum_{i=1}^{n}\sum_{i=1}^{n}\alpha_i\alpha_jy_iy_jx_ix_j + \sum_{i=1}^{n}\alpha_i\right]$$

$$\Downarrow$$

$$\min\limits_{\alpha}\left[\frac{1}{2}\sum_{i=1}^{n}\sum_{j=1}^{n}\alpha_i\alpha_jy_iy_jx_ix_j - \sum_{i=1}^{n}\alpha_i\right] \quad s.t.\begin{cases}\sum\limits_{i=1}^{n}\alpha_iy_i=0\\ c-\alpha_i-\mu_i=0\\ \alpha_i\geq 0\\ \mu_i\geq 0\end{cases}$$

整理一下得，对偶问题为：

$$\min\limits_{\alpha}\ \frac{1}{2}\sum_{i=1}^{n}\sum_{j=1}^{n}\alpha_i\alpha_jy_iy_jx_ix_j - \sum_{i=1}^{n}\alpha_i \quad s.t.\ \sum_{i=1}^{n}\alpha_iy_i=0,\ 0\leq\alpha_i\leq C$$

二、垃圾邮件分类

1、代码解析

（1）加载训练集和测试集数据

```
1.  def load_data():
2.      # 加载 mat 格式的字典文件
3.      spam_train = loadmat(file_name="spamTrain.mat")
4.      print(spam_train.keys())
5.      spam_train_x = spam_train["X"]
6.      spam_train_y = spam_train["y"]
7.      # 一个数据的长度是 1899，也就是说垃圾邮件一共有 1899 个特征
8.      spam_train_y = [math.pow(-1, i+1) for i in spam_train_y]
9.      spam_train_y = np.array(spam_train_y, dtype=int).reshape(-1, 1)
10.     # print(spam_train_y)
11.     # 同样的方式对测试集进行处理
12.     spam_test = loadmat(file_name="spamTest.mat")
13.     print(spam_test.keys())
14.     spam_test_x = spam_test["Xtest"]
15.     spam_test_y = spam_test["ytest"]
16.     spam_test_y = [math.pow(-1, i + 1) for i in spam_test_y]
17.     spam_test_y = np.array(spam_test_y, dtype=int).reshape(-1, 1)
18.     for x in spam_train_x:
19.         print("训练集特征长度:{}".format(len(x)))
20.         break
21.     for x in spam_test_x:
22.         print("测试集特征长度:{}".format(len(x)))
23.         break
24.     print("训练集样本数量:{}".format(spam_train_y.shape[0]))
25.     print("测试集样本数量:{}".format(spam_test_y.shape[0]))
26.     return spam_train_x, spam_train_y, spam_test_x, spam_test_y
```

（2）批量 Pegasos 算法，参数分别是数据，数据标签，C=0.1，训练轮数，batch 大小

```
1.  # 批量 Pegasos 算法，参数分别是数据，数据标签，C=0.1，训练轮数，batch 大小
2.  def batchPegasos(x, y, C, T, k):
3.      lam = 1 / (k * C)
4.      m, n = np.shape(x)
5.      w = np.zeros(n)
6.      dataIndex = np.array([i for i in range(m)])
7.      for t in range(1, T + 1):
8.          wDelta = np.zeros(n)
9.          eta = 1.0 / (lam * t)
10.         np.random.shuffle(dataIndex)
11.         for j in range(k):
```

```
12.             i = dataIndex[j]
13.             p = predict(w, x[i, :])
14.             if y[i][0] * p < 1:
15.                 wDelta += y[i] * x[i, :]
16.         w = (1.0 - 1 / t) * w + (eta / k) * wDelta
17.     return w
18.
19.
20. # 预测 wx+b
21. def predict(w, x):
22.     return w.T @ x
```

（3）对测试集进行测试

```
1.  def test(x, y, w):
2.   predict_y = []
3.   label_y = y.reshape(-1)
4.   # print(label_y)
5.   for x_i, y_i in zip(x, label_y):
6.       tmp = predict(w, x_i)
7.       if tmp <= 0:
8.           predict_y.append(-1)
9.       else:
10.          predict_y.append(1)
11.  predict_y = np.asarray(predict_y)
12.  # print(np.sum(predict_y == label_y))
13.  print("正确率为
     {}/{}".format(np.sum(predict_y == label_y), len(predict_y)))
```

（4）主函数

```
1.  if __name__ == '__main__':
2.   spam_train_x, spam_train_y, spam_test_x, spam_test_y = load_data()
3.   # 训练
4.   c = 0.1
5.   epochs = 100
6.   batch_size = 100
7.   w = batchPegasos(spam_train_x, spam_train_y, c, epochs, batch_size)
8.   # 测试
9.   test(spam_test_x, spam_test_y, w)
```

## 2、实验结果

```
D:\Python\anaconda\anaconda3\envs\pytorch_gpu\python.exe D:/Python/pycharm/pythonProject/softSVM/main.py
dict_keys(['__header__', '__version__', '__globals__', 'X', 'y'])
dict_keys(['__header__', '__version__', '__globals__', 'Xtest', 'ytest'])
训练集特征长度:1899
测试集特征长度:1899
训练集样本数量:4000
测试集样本数量:1000
正确率为975/1000

Process finished with exit code 0
```