

第三题

3. 试采用核映射SVM方法对“癌症患者生存期”数据集进行分类。给出具体代码，以及训练和测试的准确率结果。

```
In [1]: 1 # 数据读入
2 import numpy as np
3 import matplotlib.pyplot as plt
4 from mpl_toolkits.mplot3d import Axes3D
5
6 import matplotlib.pyplot as plt
7 import numpy as np
8 import torch
9 from sklearn import preprocessing
10 with open('haberman.data') as f:
11     raw_data = [line.strip().split(',') for line in f.readlines()]
12 np.random.seed(666)
13 raw_data = np.array(raw_data, dtype=np.int)
14 N_SAMPLES = raw_data.shape[0]
15 indices = np.arange(N_SAMPLES)
16 np.random.shuffle(indices)
17 raw_data = raw_data[indices]
18 raw_data
executed in 1.12s, finished 11:07:47 2019-12-29
```

```
Out[1]: array([[47, 66, 0, 1],
 [69, 65, 0, 1],
 [36, 69, 0, 1],
 ...,
 [46, 62, 5, 2],
 [43, 63, 2, 1],
 [61, 68, 0, 1]])
```

解决方法

线性支持向量机问题的对偶表达式为：原始优化问题转化为参数 α 的极值问题

$$A(\alpha) = \sum_{i=1}^n \alpha_i - \frac{1}{2} \sum_{i=1}^n \sum_{j=1}^n \alpha_i \alpha_j y_i y_j X_i X_j \quad (1)$$

$$st. \sum_{i=1}^n \alpha_i y_i = 0 (\alpha_i \geq 0, i = 1, 2, n) \quad (2)$$

$$y_i \in \{+1, -1\} \quad (3)$$

使用梯度的优化方法，直接最大化对偶优化问题 $\max A(\alpha)$

在前述方程里的核函数是点积， $x_i \cdot x_j$ ，其为线性核函数。该核函数是以数据点 (i, j) 的点积填充的方阵。

替代数据点间的点积，可以将其扩展到更复杂的函数更高维度。这看似不怎么复杂，但是如果选择函数 k ，其需满足如下条件：

$$k(x_i, x_j) = \varphi(x_i) \cdot \varphi(x_j) \quad (4)$$

高斯核

$$k(x_i, x_j) = e^{-\gamma \|x_i - x_j\|^2} \quad (5)$$

或者参考论文 Pegasos primal estimated sub-gradient solver 的优化公式也是可以的

$$\operatorname{argmin} \frac{\lambda}{2} \sum_{i,j} \alpha_i \alpha_j K(x_i, x_j) + \frac{1}{N} \sum_i \max\{0, 1 - y_i \sum_j \alpha_j K(x_i, x_j)\}, \alpha \geq 0 \quad (6)$$

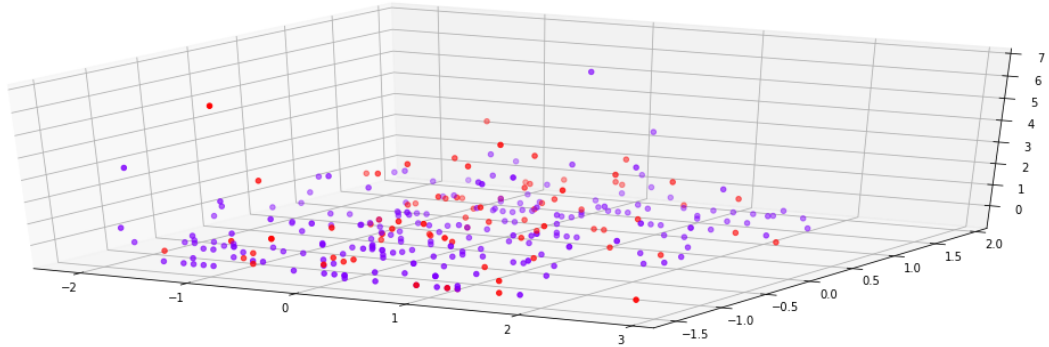
预测函数

$$\text{predict} : y = \sum \alpha_i y_i K(x, x_i) \quad (7)$$

```
In [2]: 1 x, y = raw_data[:, :-1], raw_data[:, -1] # 数据归一化
2 x = preprocessing.scale(x)
3 print(x.mean(0), x.std(0))
4
executed in 8ms, finished 11:07:47 2019-12-29
[1.43313103e-16 3.84587061e-16 1.11747938e-16] [1. 1. 1.]
```

```
In [3]: 1 fig = plt.figure(figsize=[20,6])
2 # 创建3d图形的两种方式
3 # ax = Axes3D(fig)
4 ax = fig.add_subplot(111, projection='3d')
5
6 ax.scatter(x[:,0], x[:,1], x[:,2], c=y, cmap="rainbow")
7 plt.show()
```

executed in 177ms, finished 11:07:47 2019-12-29



可见数据线性不可分，所以采用高斯核

```
In [4]: 1 y[y == 1] = -1
2 y[y == 2] = 1
3 y = y.reshape(-1, 1)
4 train_test_ratio = 0.6 # 划分训练集验证集合
5
6 N_TRAIN = round(N_SAMPLES*train_test_ratio)
7 x_train = x[:N_TRAIN]
8 y_train = y[:N_TRAIN]
9
10 x_test = x[-N_TRAIN:]
11 y_test = y[-N_TRAIN:]
```

executed in 5ms, finished 11:07:47 2019-12-29

```
In [5]: 1 def pairwise_mse(u: torch.Tensor, v: torch.Tensor):
2     ru = (u*u).sum(-1, keepdims=True)
3     rv = (v*v).sum(-1, keepdims=True)
4     return ru-2*u.mm(v.t())+rv.t()
5
6 def Linear_kernal(u, v):
7     return u.mm(v.t())
8
9 def Gaussian_kernal(u, v):
10    # sigma = 50.0
11    gamma = -10
12    return torch.exp(gamma*pairwise_mse(u, v))
13
14 kernel = Gaussian_kernal
15
16 a = torch.rand(size=[N_TRAIN, 1], requires_grad=True) # SVM paramter ai
17
18 def SVM_dul_loss(x, y):
19     aij = a.mm(a.t())
20     yij = y.mm(y.t())
21     xij = kernel(x, x)
22     assert aij.size() == yij.size() == xij.size()
23     loss = a.sum()-(0.5*aij*yij*xij).sum()
24     return - loss
25
26 # def SVM_dul_loss(x, y, lambda_=0.5):
27 #     aij = a.mm(a.t())
28 #     yij = y.mm(y.t())
29 #     xij = kernel(x, x)
30 #     assert aij.size() == yij.size() == xij.size()
31 #     loss = (lambda_ / 2 * aij * xij).sum() + \
32 #         (1-xij.mm(a)*y).clamp(min=0).mean()
33 #     return loss
34
35 def SVM_dul_predict(x_ture, y_true, x_pred):
36     xij = kernel(x_ture, x_pred)
37     yai = a*y_true
38     y_pred = yai.t().mm(xij)
39     y_pred = torch.sign(y_pred-y_pred.mean())
40     return y_pred.view(-1, 1)
41
```

executed in 29ms, finished 11:07:47 2019-12-29

```

In [6]: 1 opt = torch.optim.SGD([a], lr=0.002)
2 x_train = torch.tensor(x_train, dtype=torch.float32)
3 y_train = torch.tensor(y_train, dtype=torch.float32)
4 x_test = torch.tensor(x_test, dtype=torch.float32)
5 y_test = torch.tensor(y_test, dtype=torch.float32)
6 loss_history = []
7 train_acc = []
8 val_acc = []
9 for epoch in range(10000):
10     loss = SVM_dul_loss(x_train, y_train)
11     loss.backward()
12     opt.step()
13     opt.zero_grad()
14     with torch.no_grad():
15         loss_history.append(loss.item())
16         y_pred = SVM_dul_predict(x_train, y_train, x_train)
17         acc = (y_train == y_pred).float().mean()
18         train_acc.append(acc)
19         # print(acc)
20         # print(a[0])
21     if epoch % 100 == 0:
22         y_pred = SVM_dul_predict(x_train, y_train, x_test)
23         acc = (y_test == y_pred).float().mean()
24         val_acc.append(acc)

```

executed in 15.0s, finished 11:08:02 2019-12-29

```

In [7]: 1 with torch.no_grad():
2     y_pred = SVM_dul_predict(x_train, y_train, x_test)
3     acc = (y_test == y_pred).float().mean()
4     print("acc", acc)
5     # Plot batch accuracy
6     fig = plt.figure(figsize=[20, 5])
7     ax = fig.add_subplot(131)
8     ax.plot(train_acc, 'k-', label='Accuracy')
9     ax.set_title('training Accuracy')
10    ax.set_xlabel('Generation')
11    ax.set_ylabel('Accuracy')
12    ax.legend(loc='lower right')
13
14    # Plot loss over time
15    ax = fig.add_subplot(132)
16    ax.plot(loss_history, 'k-')
17    ax.set_title('Loss per Generation')
18    ax.set_xlabel('Generation')
19    ax.set_ylabel('Loss')
20
21    ax = fig.add_subplot(133)
22    ax.plot(val_acc, 'k-', label='Accuracy')
23    ax.set_title('validation Accuracy')
24    ax.set_xlabel('Generation')
25    ax.set_ylabel('Accuracy')
26    ax.legend(loc='lower right')
27    plt.show()

```

executed in 398ms, finished 11:08:02 2019-12-29

acc tensor(0. 7391)

