

SVM心路历程

1. SVM原理推导

- 1. 1. 定义任务
- 1. 2. 决策函数标准化
- 1. 3. 计算街宽
- 1. 4. 定义优化目标
- 1. 5. 定义拉格朗日函数
- 1. 6. 核方法
- 1. 7. 神经网络

2. 对偶问题转化

3. SMO算法

4. 附录

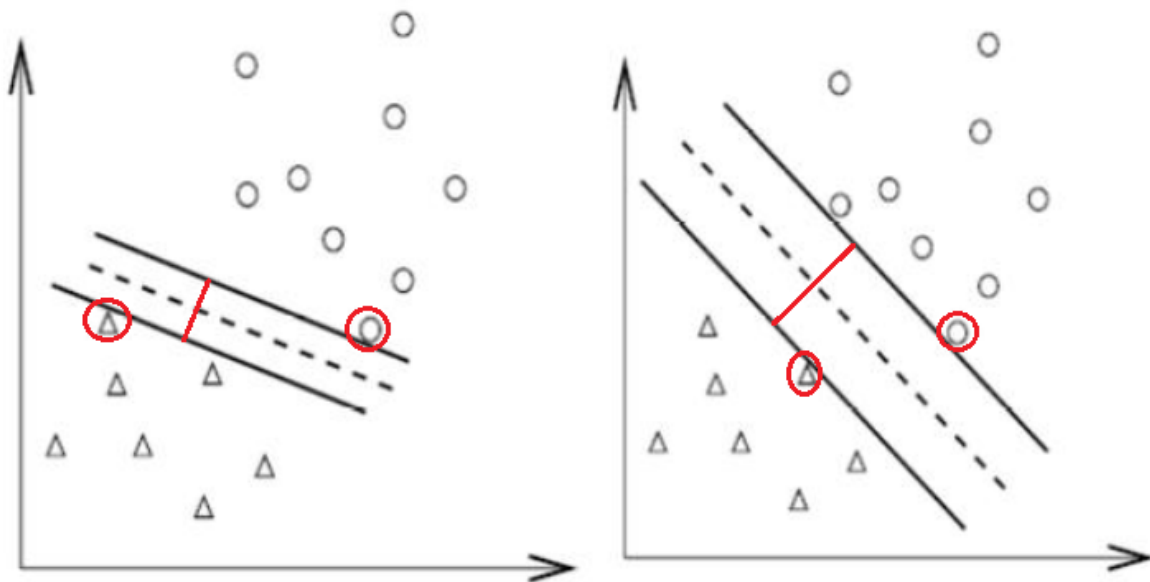
- 4. 1. 二范数求导

1. SVM原理推导

1.1. 定义任务

任务：对线性可分的二分类问题，寻找一条直线，对于距离该直线最近的正负样本点，使得它们到直线的距离相等，且它们之间的距离最大。**即寻找最宽的路，将正负样本分开。**

如下图所示，右侧的决策边界要优于左侧决策边界。



设 \vec{w} 为决策边界的法向量， C 为原点到决策边界的距离，则决策函数为：

$$\begin{cases} \vec{w} \cdot \vec{x} > C, \text{ then } + \\ \vec{w} \cdot \vec{x} < C, \text{ then } - \end{cases} \quad (27)$$

稍作化简，可得：

$$\begin{cases} \vec{w} \cdot \vec{x} + b > 0, \text{ then } + \\ \vec{w} \cdot \vec{x} + b < 0, \text{ then } - \end{cases} \quad (28)$$

对 y 做如下定义：

$$y = \begin{cases} 1, + \\ -1, - \end{cases} \quad (29)$$

则决策函数应该满足：

$$y(\vec{w} \cdot \vec{x} + b) > 0 \quad (30)$$

1.2. 决策函数标准化

对决策函数加入距离限制，比如要求最近的点到决策边界的函数间隔至少为1，则决策函数：

$$\begin{aligned} y(\vec{w} \cdot \vec{x} + b) &\geq 1 \\ x &\in \text{训练集} \end{aligned} \quad (31)$$

含义：正负样本点距离决策边界的最近距离为1。对于等于1的点，我们称之为支持向量。

1.3. 计算街宽

$$\begin{aligned} WIDTH &= (\vec{x}_+ - \vec{x}_-) \cdot \frac{\vec{w}}{|\vec{w}|} \\ &= \frac{2}{|\vec{w}|} \quad \vec{x}_+ \text{ 和 } \vec{x}_- \text{ 均为支持向量} \end{aligned} \quad (32)$$

1.4. 定义优化目标

$$\begin{aligned} & \max \frac{2}{|\vec{w}|} \\ & \Rightarrow \min |\vec{w}| \\ & \Rightarrow \begin{cases} \min \frac{1}{2} \|\vec{w}\|_2^2 \\ s. t. \quad y_i(\vec{w} \cdot \vec{x}_i + b) \geq 1, i = 1 \dots N \end{cases} \end{aligned} \quad (33)$$

1.5. 定义拉格朗日函数

$$\begin{aligned} L(\alpha, w, b) &= \frac{1}{2} \|\vec{w}\|_2^2 + \sum_{i=1}^N \alpha_i [1 - y_i(\vec{w} \cdot \vec{x}_i + b)] \\ \vec{\alpha} &= (\alpha_1, \alpha_2, \dots, \alpha_N) \geq 0 \end{aligned} \quad (34)$$

优化目标是使 L 最小，使用数值方法， L 分别对 \vec{w} 和 b 求导：

$$\begin{aligned} \frac{\partial L}{\partial \vec{w}} &= \vec{w} - \sum_{i=1}^N \alpha_i y_i x_i = 0 \\ \frac{\partial L}{\partial b} &= - \sum_{i=1}^N \alpha_i y_i = 0 \end{aligned} \quad (35)$$

得 \vec{w} 的最优解：

$$\begin{aligned} \vec{w}^* &= \sum_{i=1}^N \alpha_i y_i x_i \\ \sum_{i=1}^N \alpha_i y_i &= 0 \end{aligned} \quad (36)$$

对 b 的最优解，可以任意找一个支持向量，根据 $y_i(\vec{w} \cdot \vec{x}_i + b) = 1$ 进行求解。

至此，求得决策边界：

$$\begin{cases} \sum_{i=1}^N \alpha_i y_i (\vec{x}_i \cdot \vec{x}_{new}) + b > 0, \text{ then } + \\ \sum_{i=1}^N \alpha_i y_i (\vec{x}_i \cdot \vec{x}_{new}) + b < 0, \text{ then } - \end{cases} \quad (37)$$

1.6. 核方法

$$L = \sum_{i=1}^N \alpha_i - \frac{1}{2} \sum_{i=1}^N \sum_{j=1}^N \alpha_i \alpha_j y_i y_j (x_i \cdot x_j), x \in \text{训练集} \quad (38)$$

注：由此可见， L 取决于训练集 $x_i \cdot x_j$

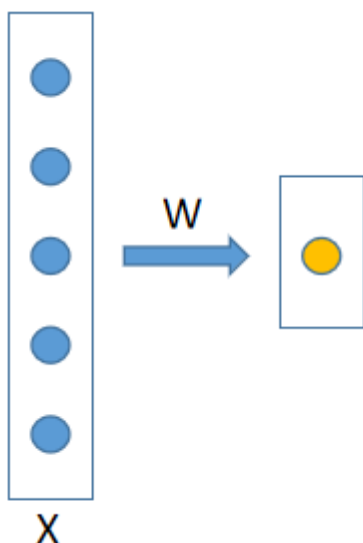
定义：

$$K(x_i, x_j) = \Phi(x_i) \cdot \Phi(x_j) \quad (39)$$

则 $K(x_i, x_j)$ 为核函数， Φ 为变换函数。左侧为点乘后的变换，右侧为变换后的点乘，核方法本质是将变换后的点乘改为点乘后的变换。我们不需要知道如何变换，只需要知道变换后的结果。

1.7. 神经网络

从神经网络的角度理解，SVM的本质是通过全连接变换，采用hinge-loss作为损失函数的分类问题。



hinge-loss（经验风险）为：

$$L(y \cdot (w \cdot x + b)) = [1 - y \cdot (w \cdot x + b)]_+ \quad (40)$$

其中， $[z]_+ = \begin{cases} z, & z > 0 \\ 0, & z \leq 0 \end{cases}$

为了防止过拟合，实际优化目标还会加入正则化项（结构风险）：

$$\min \{ [1 - y \cdot (w \cdot x + b)]_+ + \lambda \|w\|^2 \}_{w,b} \quad (41)$$

其中， $[z]_+ = \begin{cases} z, & z > 0 \\ 0, & z \leq 0 \end{cases}$

实现代码为：

```
1 class LinearSVM(nn.Module):
2     def __init__(self):
3         super(LinearSVM, self).__init__()
4         self.linear = nn.Linear(in_features=2, out_features=1)
5     def forward(self, x):
6         y = self.linear(x)
7         return y
8
9 svm = LinearSVM()
10 optimizer = optim.SGD(svm.parameters(), lr=0.1)
11
12 batch_size = 1
13 epoch_num = 30
14 N = 500
15
16 for epoch in range(1, epoch_num+1):
17     for i in range(0, N, batch_size):
18         x = torch.Tensor(X[i:i+batch_size])
19         y = torch.Tensor(Y[i:i+batch_size])
20
```

```

21     y_pred = svm(x)
22
23     loss = torch.mean(torch.clamp(1 - y_pred * y, min=0)) # hinge
24     loss += 0.01 * torch.mean(svm.linear.weight ** 2) / 2 # l2
25
26     loss.backward()
27     optimizer.step()
28     optimizer.zero_grad()
29

```

2. 对偶问题转化

1.原始问题:

$$L(\alpha, w, b) = \frac{1}{2} \|\vec{w}\|_2^2 + \sum_{i=1}^N \alpha_i [1 - y_i (\vec{w} \cdot \vec{x}_i + b)] \quad (42)$$

$$\vec{\alpha} = (\alpha_1, \alpha_2, \dots, \alpha_N) \geq 0$$

2.对偶函数:

$$g(\alpha) = \min_{w, b} L(\alpha, w, b) \quad (43)$$

3.对偶问题:

$$\max_{\alpha, \alpha \geq 0} g(\alpha) = \max_{\alpha, \alpha \geq 0} \min_{w, b} L(\alpha, w, b) \quad (44)$$

4.求解对偶函数:

$$\nabla_w L(w, b, \alpha) = w - \sum_{i=1}^N \alpha_i y_i x_i = 0 \quad (45)$$

$$\nabla_b L(w, b, \alpha) = - \sum_{i=1}^N \alpha_i y_i = 0$$

解得:

$$w = \sum_{i=1}^N \alpha_i y_i x_i \quad (46)$$

$$\sum_{i=1}^N \alpha_i y_i = 0$$

对 $L(\alpha, w, b)$ 化简得:

$$L(w, b, \alpha) = \frac{1}{2} \sum_{i=1}^N \sum_{j=1}^N \alpha_i \alpha_j y_i y_j (x_i \cdot x_j) - \sum_{i=1}^N \alpha_i y_i \left(\left(\sum_{j=1}^N \alpha_j y_j x_j \right) \cdot x_i + b \right) + \sum_{i=1}^N \alpha_i \quad (47)$$

$$= -\frac{1}{2} \sum_{i=1}^N \sum_{j=1}^N \alpha_i \alpha_j y_i y_j (x_i \cdot x_j) + \sum_{i=1}^N \alpha_i$$

所以, 对偶函数为:

$$g(\alpha) = \min_{w,b} L(\alpha, w, b) = -\frac{1}{2} \sum_{i=1} \sum_{j=1} \alpha_i \alpha_j y_i y_j (x_i \cdot x_j) + \sum_{i=1} \alpha_i \quad (48)$$

5.对偶函数极大化，即将原问题转化为对偶问题：

$$\begin{aligned} \max_{\alpha} \quad & -\frac{1}{2} \sum_{i=1}^N \sum_{j=1}^N \alpha_i \alpha_j y_i y_j (x_i \cdot x_j) + \sum_{i=1}^N \alpha_i \\ \text{s.t.} \quad & \sum_{i=1}^N \alpha_i y_i = 0 \\ & \alpha_i \geq 0, \quad i = 1, 2, \dots, N \end{aligned} \quad (49)$$

将问题改为等价的 \min 问题：

$$\begin{aligned} \min_{\alpha} \quad & \frac{1}{2} \sum_{i=1}^N \sum_{j=1}^N \alpha_i \alpha_j y_i y_j (x_i \cdot x_j) - \sum_{i=1}^N \alpha_i \\ \text{s.t.} \quad & \sum_{i=1}^N \alpha_i y_i = 0 \\ & \alpha_i \geq 0, \quad i = 1, 2, \dots, N \end{aligned} \quad (50)$$

3. SMO算法

待补充。

4. 附录

4.1. 二范数求导

$$\|\vec{w}\|_2^2 = \vec{w}^T \cdot \vec{w} = w_1^2 + w_2^2 + \dots + w_n^2 \quad (51)$$

$$\frac{\partial \vec{w}^T \cdot \vec{w}}{\partial \vec{w}} = 2\vec{w} \quad (52)$$

解释： $w_1^2 + w_2^2 + \dots + w_n^2$ 分别对 \vec{w} 中的每个变量求偏导。