

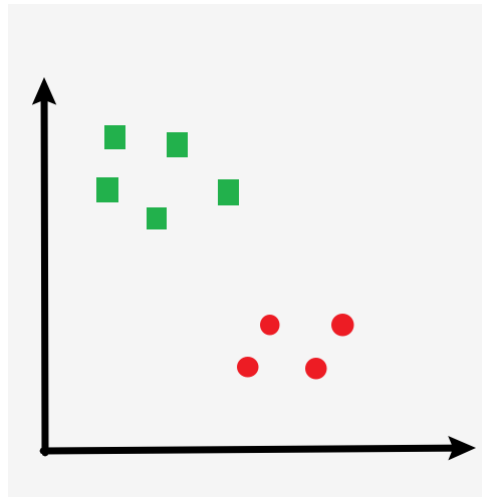
# Support Vector Machine

Baicheng Chen

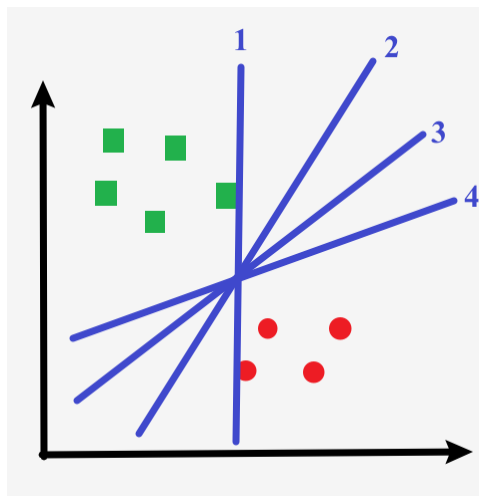
支持向量机是一个经典的二分类模型，由Cortes和Vapnik于1995年提出。它区别于感知机，感知机是通过错误驱动的方式来确定一个可行的决策边界（**minimize loss function** => **misclassification points**，有无穷多个），而支持向量机则是选出“最优”的一个决策边界。这里的“最优”如何定义呢？SVM考虑了最大间隔（两个类别之间）。

这里，我们首先考虑一个简单的线性可分的二分类问题。

假设有样本点  $\{(x_i, y_i)\}_i^N, x_i \in R^p, y \in \{-1, 1\}$ 。



现在，若我们将绿色样本点与红色样本点分开，有无数条决策边界可以做到。下图中的1、2、3、4这四条线都可以完美的将这两种样本点进行分类。但我们可以发现1号线明显鲁棒性较差，如果样本点变多，它的泛化误差较大，并不能算一个足够“好”的边界。而我们需要边界是一条鲁棒性好，可泛化的决策边界。



于是，SVM提出了“最大间隔”的概念。从上图中看，最大间隔便是两个类别最“中间”的那条线（即为3号线），它保持了到两个分类区域的距离最大，我们认为“最大间隔”的决策边界就是“最好”的决策边界。而这就引出了这篇文章的第一部分，也就是线性可分支持向量机，也称之为最大间隔分类器（Hard-margin SVM）。

## Hard-Margin SVM

现在，我们从数学的角度来定义“最大间隔”这个概念。首先，我们可以将决策边界看作一个超平面，用公式表示为  $0 = w^T x + b$ ，这个超平面将  $n$  维空间分割为两半，其中，法向量  $w$  指向的那一半定义为正空间 ( $\forall x^+, w^T x^+ + b > 0$ )，反之另一半则为负空间，可知：

$$\begin{aligned} & \max \text{margin}(w, b) \\ & s. t. \begin{cases} w^T x_i + b > 0 & , y_i = +1 \\ w^T x_i + b < 0 & , y_i = -1 \end{cases} \end{aligned}$$

这个公式还是不够简洁，我们尝试将 `constraints` 合并为一个条件，即：

$$\begin{aligned} & \max \text{margin}(w, b) \\ & s. t. \quad y_i(w^T x_i + b) > 0, \text{ for } \forall i = 1, 2, \dots, N \end{aligned}$$

这个条件与上面两个条件是等价的，因为同号相乘一定为正 ( $y_i$  与  $w^T x_i + b$  始终同号)。

接下来，我们来详细推导  $\text{margin}(w, b)$  这个函数，即解释如何用数学形式来表示我们在前文所提到的“间隔”。简单来讲，“间隔”就是分类区域到决策边界的距离，在这里，我们认为  $N$  个样本点到决策边界的距离中最小的那个即为“间隔”。这样，如果我们通过最大化这个最小距离找到了一个最优的决策边界，它对于整个分类区域来说也一定是最优的。

假设这个最小距离由样本点  $(x_i, y_i)$  提供，下面计算  $(x_i, y_i)$  到决策边界  $y = w^T x + b$  的距离：

对于点  $x_0$ ，设其在超平面  $0 = w^T x + b$  的投影为  $x_1$ ，则有  $w^T x_1 + b = 0$ 。

因为法向量  $w$  垂直于超平面， $x_1 - x_0$  也垂直于超平面，故  $w \parallel x_1 - x_0$ ，则：

$$\begin{aligned} |w \cdot x_1 - x_0| &= ||w|| \cdot \cos \pi \cdot ||x_1 - x_0|| = ||w|| \cdot ||x_1 - x_0|| = ||w|| \cdot r, \quad r \text{ 为距离;} \\ w \cdot x_1 - x_0 &= w_1(x_1^0 - x_1^1) + w_2(x_2^0 - x_2^1) + \dots + w_n(x_n^0 - x_n^1) \\ &= w_1 x_1^0 + w_2 x_2^0 + \dots + w_n x_n^0 - (w_1 x_1^1 + w_2 x_2^1 + \dots + w_n x_n^1) \\ &= w^T x_0 - w^T x_1 \\ &= w^T x_0 + b \quad (\because w^T x_1 + b = 0). \end{aligned}$$

所以有，

$$\begin{aligned} |w^T x_0 + b| &= |w \cdot x_1 - x_0| \\ &= ||w|| \cdot r \\ \Rightarrow r &= \frac{|w^T x_i + b|}{||w||} \end{aligned}$$

有了这个距离，我们就可以表示  $\text{margin}(w, b)$  函数了，注意： $\text{margin}$  是样本点到决策边界的最小距离！

$$\begin{aligned} \text{margin}(w, b) &= \min_{w, b, x_i} \text{distance}(w, b, x_i) \\ &= \min_{w, b, x_i} \frac{|w^T x_i + b|}{||w||} \end{aligned}$$

将导出的  $\text{margin}(w, b)$  公式带回到我们最开始写出的优化问题中：

$$\begin{aligned} & \max_{w, b} \min_{x_i} \frac{|w^T x_i + b|}{||w||} \\ & s. t. \quad y_i(w^T x_i + b) > 0 \end{aligned}$$

因为我们有 `constrain`  $y_i(w^T x_i + b) > 0$ ，则可以将  $y_i$  看作  $w^T x_i + b$  的绝对值符号，将其替换到 `objective function` 中：

$$\begin{aligned} \max_{w,b} \min_{x_i} \frac{y_i(w^T x_i + b)}{\|w\|} \\ s.t. \quad y_i(w^T x_i + b) > 0 \end{aligned}$$

因为 $y_i(w^T x_i + b)$ 始终大于0，我们假设 $\exists \gamma > 0, s.t. \min_{x_i, y_i} y_i(w^T x_i + b) = \gamma$ 。

则我们可以将objective function转化为：

$$\begin{aligned} \max_{w,b} \min_{x_i} \frac{y_i(w^T x_i + b)}{\|w\|} &= \max_{w,b} \frac{1}{\|w\|} \min_{x_i} y_i(w^T x_i + b) \\ &= \max_{w,b} \frac{1}{\|w\|} \gamma \end{aligned}$$

我们知道，超平面方程不唯一，即当我们等倍缩放 $w$ 和 $b$ 时，所得的新超平面与原超平面相同，故而， $\gamma$ 有无数种可能的取值。所以，我们添加一个约束，即令 $\gamma = 1$ ，以此来限制到唯一超平面。优化问题变成了：

$$\begin{aligned} \max_{w,b} \frac{1}{\|w\|} \\ s.t. \quad \min y_i(w^T x_i + b) = 1 \end{aligned}$$

等价于，

$$\begin{aligned} \max_{w,b} \frac{1}{\|w\|} \\ s.t. \quad y_i(w^T x_i + b) \geq 1 \end{aligned}$$

等价于，

$$\begin{aligned} \min_{w,b} \|w\| \\ s.t. \quad y_i(w^T x_i + b) \geq 1 \end{aligned}$$

显然，这个优化问题可以被转化成一个QP问题，即：

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

等价于，

$$\begin{aligned} \min_{w,b} \frac{1}{2} w^T w \\ s.t. \quad 1 - y_i(w^T x_i + b) \leq 0, \text{ for } \forall i = 1, \dots, N \end{aligned}$$

将这个QP问题看作一个约束优化问题，而上式则为该约束优化问题的原问题 (primal problem)。

下面，我们引入拉格朗日函数，

$$L(w, b, \lambda) = \frac{1}{2} w^T w + \sum_{i=1}^N \lambda_i (1 - y_i(w^T x_i + b)), \quad \lambda_i \geq 0$$

通过拉格朗日函数的引入，我们可以找出原问题的无约束形式：

$$\begin{aligned} \min_{w,b} \max_{\lambda} L(w, b, \lambda) \\ s.t. \quad \lambda_i \geq 0 \end{aligned}$$

下面对上述优化问题中的objective function进行解释：

$$\text{对于 } L(w, b, \lambda) = \frac{1}{2} w^T w + \sum_{i=1}^N \lambda_i (1 - y_i (w^T x_i + b)),$$

我们在左右两边同时取max, 即：

$$\max_{\lambda} L(w, b, \lambda) = \frac{1}{2} w^T w + \max_{\lambda} \left( \sum_{i=1}^N \lambda_i (1 - y_i (w^T x_i + b)) \right)$$

在该式中,  $\lambda_i \geq 0, 1 - y_i (w^T x_i + b) \leq 0$ , 我们可以得到,

$$\lambda_i (1 - y_i (w^T x_i + b)) \leq 0$$

$$\text{即 } \max_{\lambda} \left( \sum_{i=1}^N \lambda_i (1 - y_i (w^T x_i + b)) \right) = 0$$

由此我们可以推出,

$$\begin{aligned} \frac{1}{2} w^T w &= \max_{\lambda} L(w, b, \lambda) - \max_{\lambda} \left( \sum_{i=1}^N \lambda_i (1 - y_i (w^T x_i + b)) \right) \\ &= \max_{\lambda} L(w, b, \lambda) - 0 \\ &= \max_{\lambda} L(w, b, \lambda) \end{aligned}$$

于是, 我们就可以得出原问题的对偶问题 (dual problem)。因为该优化问题为凸二次优化问题, 且约束条件为仿射函数, 满足放松Slater条件, 故原问题与对偶问题满足强对偶关系 (convex + slater  $\Rightarrow$  strong duality)。

$$\begin{aligned} \max_{\lambda} \min_{w, b} L(w, b, \lambda) \\ \text{s.t. } \lambda_i \geq 0 \end{aligned}$$

下面, 我们着手解决这个无约束优化问题。

- 对b求导,

$$\begin{aligned} \frac{\partial L}{\partial b} &= \frac{\partial}{\partial b} \left[ \sum_{i=1}^N \lambda_i - \sum_{i=1}^N \lambda_i y_i (w^T x_i + b) \right] \\ &= \frac{\partial}{\partial b} \left[ - \sum_{i=1}^N \lambda_i y_i b \right] \\ &= - \sum_{i=1}^N \lambda_i y_i = 0 \end{aligned}$$

- 将上述结果带入原式化简,

$$\begin{aligned} L(w, b, \lambda) &= \frac{1}{2} w^T w + \sum_{i=1}^N \lambda_i (1 - y_i (w^T x_i + b)) \\ &= \frac{1}{2} w^T w + \sum_{i=1}^N \lambda_i - \sum_{i=1}^N \lambda_i y_i (w^T x_i + b) \\ &= \frac{1}{2} w^T w + \sum_{i=1}^N \lambda_i - \sum_{i=1}^N \lambda_i y_i w^T x_i - \sum_{i=1}^N \lambda_i y_i b \\ &= \frac{1}{2} w^T w + \sum_{i=1}^N \lambda_i - \sum_{i=1}^N \lambda_i y_i w^T x_i \quad (\because - \sum_{i=1}^N \lambda_i y_i = 0 \text{ 且 } b \text{ 为常数}) \end{aligned}$$

- 对w求导,

$$\frac{\partial L}{\partial w} = \frac{1}{2} \cdot 2 \cdot w - \sum_{i=1}^N \lambda_i y_i x_i = 0$$

$$\Rightarrow w^* = \sum_{i=1}^N \lambda_i y_i x_i$$

- 将 $w^*$ 代入objective function,

$$\begin{aligned} \min_{w,b} L(w, b, \lambda) &= \frac{1}{2} \left( \sum_{i=1}^N \lambda_i y_i x_i \right)^T \left( \sum_{j=1}^N \lambda_j y_j x_j \right) - \sum_{i=1}^N \lambda_i y_i \left( \sum_{j=1}^N \lambda_j y_j x_j \right)^T x_i + \sum_{i=1}^N \lambda_i \\ &= \frac{1}{2} \sum_{i=1}^N \sum_{j=1}^N \lambda_i \lambda_j y_i y_j x_i^T x_j - \sum_{i=1}^N \sum_{j=1}^N \lambda_i \lambda_j y_i y_j x_j^T x_i + \sum_{i=1}^N \lambda_i \\ &= -\frac{1}{2} \sum_{i=1}^N \sum_{j=1}^N \lambda_i \lambda_j y_i y_j x_i^T x_j + \sum_{i=1}^N \lambda_i \end{aligned}$$

- 由上述推导，对偶问题可化为，

$$\begin{aligned} \max_{\lambda} \quad & -\frac{1}{2} \sum_{i=1}^N \sum_{j=1}^N \lambda_i \lambda_j y_i y_j x_i^T x_j + \sum_{i=1}^N \lambda_i \\ \text{s.t.} \quad & \lambda_i \geq 0 \end{aligned}$$

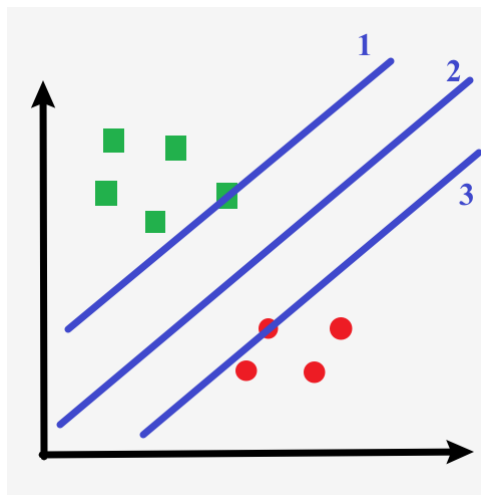
- 等价于，

$$\begin{aligned} \min_{\lambda} \quad & \frac{1}{2} \sum_{i=1}^N \sum_{j=1}^N \lambda_i \lambda_j y_i y_j x_i^T x_j + \sum_{i=1}^N \lambda_i \\ \text{s.t.} \quad & \lambda_i \geq 0 \end{aligned}$$

因为原问题与对偶问题具有强对偶关系 $\Leftrightarrow$ 满足KKT条件:

$$KKT \begin{cases} \text{可行条件: } \begin{cases} \lambda_i \geq 0 \\ 1 - y_i(w^T x + b) \leq 0 \end{cases} \\ \text{互补松弛条件: } \lambda_i (1 - y_i(w^T x + b)) = 0 \\ \text{梯度为0: } \frac{\partial L}{\partial w} = 0, \frac{\partial L}{\partial b} = 0 \end{cases}$$

现在，我们再来看图示。中间蓝线便是我们要找出的最优决策边界 $w^T x + b = 0$ （线2），根据前文所述，我们将 $\gamma$ 限制在1，故而，经过距离最优决策边界最近的两个样本点的直线可以分别表示为 $w^T x + b = 1$ （线1）和 $w^T x + b = -1$ （线3）。



以  $w^T x + b = 1$  为例（绿色区域），由可行条件  $1 - y_i(w^T x + b) \leq 0$  可知，此时为取等条件，即  $1 - y_i(w^T x + b) = 0$ 。由互补松弛条件可知， $\lambda$  可取任意值，满足  $\lambda_i \geq 0$ 。下面我们看到线1上方的区域，这里仍然存在四个绿色样本点，对于它们来说， $w^T x + b > 1$ ，所以  $1 - y_i(w^T x + b) < 0$ ，故而由互补松弛条件可知， $\lambda_i$  一定等于零。所以，对于  $\lambda_i$  的值有贡献的样本点仅存在于  $w^T x + b = 1$  和  $w^T x + b = -1$  两条线上，我们称这些点为支持向量。

下面，我们做最后的结果推导。根据KKT条件以及前文推导，我们已知：

$$w^* = \sum_{i=1}^N \lambda_i y_i x_i$$

从1和3两条线入手，进行  $b^*$  的推导：

我们假设  $\exists (x_k, y_k), s.t. 1 - y_k(w^T x_k + b) = 0 \Rightarrow y_k(w^T x_k + b) = 1$

因为  $y_k = \pm 1$ ，等式两边同乘  $y_k$  仍然成立，则：

$$y_k^2(w^T x_k + b) = y_k$$

$$\Rightarrow b^* = y_k - w^T x_k = y_k - \sum_{i=1}^N \lambda_i y_i x_i^T x_k$$

至此，我们便完成了Hard-margin SVM的全部推导，该判别模型最终可表示为一下分类决策函数：

$$f(x) = \text{sign}(w^{*T} x + b^*)$$

$$w^* = \sum_{i=1}^N \lambda_i y_i x_i$$

$$b^* = y_k - w^T x_k = y_k - \sum_{i=1}^N \lambda_i y_i x_i^T x_k$$

## Soft-Margin SVM

在引入部分我们便给Hard-margin SVM添加了一个限制条件，即我们的二分类问题一定是线性可分的。那么如果我们拿到的二分类问题是线性不可分问题呢？这里，我们就引入了Soft-margin SVM，它相当于在线性可分问题的基础上允许一些错误样本点（整体分类问题是线性不可分的）。

回顾我们在推导Hard-margin SVM时的原问题，即：

$$\begin{aligned} \min_{w,b} \quad & \frac{1}{2} w^T w \\ s.t. \quad & 1 - y_i(w^T x_i + b) \leq 0, \text{ for } \forall i = 1, \dots, N \end{aligned}$$

Soft-margin SVM的基本思想是引入一个loss function用来表示错误分类损失，并将其一起最小化，获得尽可能优的决策边界，我们可以写作：

$$\begin{aligned} \min_{w,b} \quad & \frac{1}{2} w^T w + \text{loss function} \\ s.t. \quad & 1 - y_i(w^T x_i + b) \leq 0, \text{ for } \forall i = 1, \dots, N \end{aligned}$$

下面我们来探讨合适的loss function。

首先，最为简单的loss function便是错分的样本点数量，数学公式可以写作：

$$\text{loss} = \sum_{i=1}^N I\{y_i(w^T x_i + b) < 1\}$$

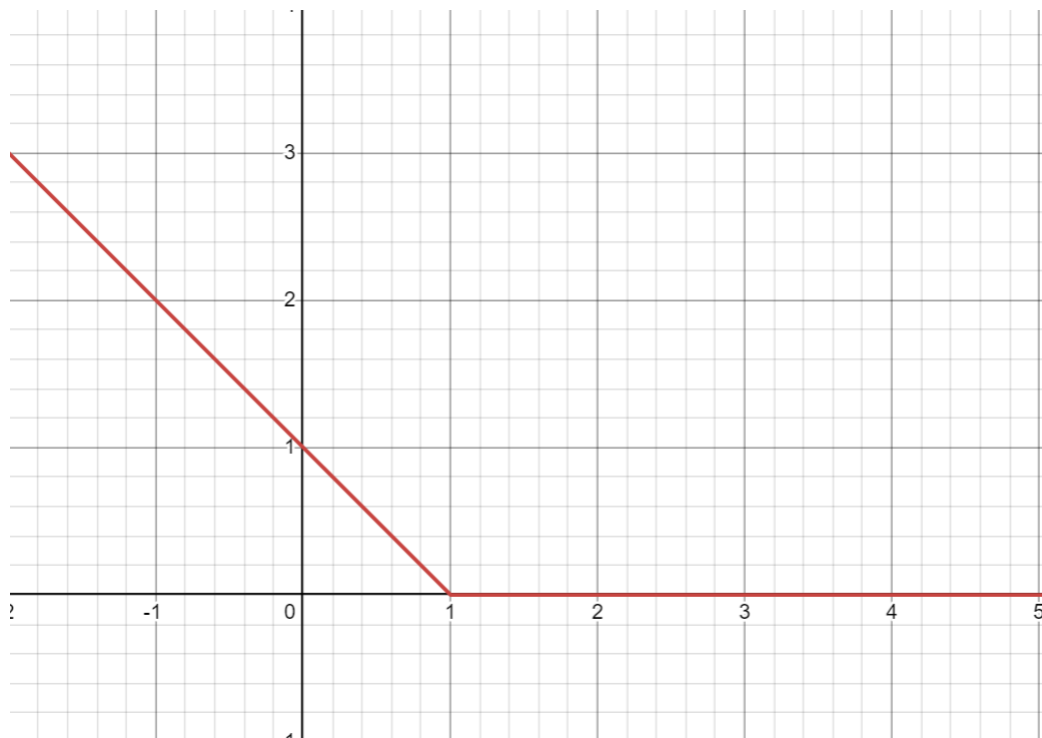
但此时的问题是这个函数并不连续，这样的优化问题很难求解，我们选择考虑别的loss function。

既然数量不行，距离行不行呢？这里，我们便引入了Hinge Loss。我们做如下考虑，假设我们已分类了样本点  $(x_i, y_i)$ ，如果该样本点满足约束条件（函数间隔（确信度）），即  $y_i(w^T x_i + b) \geq 1$ ，令  $loss = 0$ ；相反，如果该样本点不满足约束条件，即  $y_i(w^T x_i + b) < 1$ ，令  $loss = 1 - y_i(w^T x_i + b)$ 。

合并成一个连续的loss function，

$$loss = \max\{0, 1 - y_i(w^T x_i + b)\}$$

设  $z = y_i(w^T x_i + b)$ ， $loss = \max\{0, 1 - z\}$ ，画出图像如下（合页形状）：



我们发现，这个损失函数是连续可微的，可以使用，于是写出Soft-margin SVM的优化问题如下：

$$\begin{aligned} \min_{w,b} \quad & \frac{1}{2} w^T w + C \sum_{i=1}^N \max\{0, 1 - y_i(w^T x_i + b)\} \\ \text{s.t.} \quad & y_i(w^T x_i + b) \geq 1, \text{ for } \forall i = 1, \dots, N \end{aligned}$$

这里， $C > 0$ 是一个超参数，我们称为惩罚参数，由实际问题决定。C越大，对错误分离的惩罚越大。

为进一步消除max函数，我们引入一个松弛变量  $\xi_i = 1 - y_i(w^T x_i + b)$ ， $\xi_i \geq 0$ 。上式可简化为，

$$\begin{aligned} \min_{w,b} \quad & \frac{1}{2} w^T w + C \sum_{i=1}^N \xi_i \\ \text{s.t.} \quad & y_i(w^T x_i + b) \geq 1, \text{ for } \forall i = 1, \dots, N \\ & \xi_i \geq 0 \end{aligned}$$

接下来，使用约束优化问题的方法同Hard-margin SVM一样正常推导即可。

## Kernel Method

前文中，我们讨论了对于线性可分二分类问题、线性不可分二分类问题的解决方法，现在，如果二分类问题变成了非线性问题，我们该怎么解决呢？对于这种问题，我们便要引入核方法。因篇幅有限，我会在下一篇文章详细讨论核方法的原理、核函数及正定核等内容。