

# SVM

## 1. 基本概念

支持向量机 (Support Vector Machine, SVM) 的基本模型是在特征空间上找到最佳的分离超平面使得训练集上正负样本间隔最大。SVM是用来解决二分类问题的有监督学习算法，在引入了核方法之后SVM也可以用来解决非线性问题。

一般SVM有下面三种：

- 硬间隔支持向量机（线性可分支持向量机）：当训练数据线性可分时，可通过硬间隔最大化求得一个线性可分支持向量机。
- 软间隔支持向量机：当训练数据近似线性可分时，可通过软间隔最大化求得一个线性支持向量机。
- 非线性支持向量机：当训练数据线性不可分时，可通过核方法以及软间隔最大化求得一个非线性支持向量机。

## 2. 硬间隔支持向量机

给定训练样本集  $D = \{(\vec{x}_1, y_1), (\vec{x}_2, y_2), \dots, (\vec{x}_n, y_n)\}$   $D = \{(x_1 \rightarrow, y_1), (x_2 \rightarrow, y_2), \dots, (x_n \rightarrow, y_n)\}$ ,  $y_i \in \{+1, -1\}$   $y_i \in \{+1, -1\}$ ,  $i$  表示第  $i$  个样本， $n$  表示样本容量。分类学习最基本的想法就是基于训练集  $D$  在特征空间中找到一个最佳划分超平面将正负样本分开，而SVM算法解决的就是如何找到最佳超平面的问题。超平面可通过如下的线性方程来描述：

$$\vec{w}^T \vec{x} + b = 0 \quad (1)$$

$$(1) \vec{w} \rightarrow T \vec{x} \rightarrow + b = 0$$

其中  $\vec{w}$  表示法向量，决定了超平面的方向； $b$  表示偏移量，决定了超平面与原点之间的距离。

对于训练数据集  $D$  假设找到了最佳超平面  $\vec{w}^* \vec{x} + b^* = 0$   $\vec{w}^* \rightarrow \vec{x} \rightarrow + b^* = 0$ ，定义决策分类函数

$$f(\vec{x}) = \text{sign}(\vec{w}^* \vec{x} + b^*) \quad (2)$$

$$(2) f(\vec{x} \rightarrow) = \text{sign}(\vec{w}^* \rightarrow \vec{x} \rightarrow + b^*)$$

该分类决策函数也称为线性可分支持向量机。

在测试时对于线性可分支持向量机可以用一个样本离划分超平面的距离来表示分类预测的可靠程度，如果样本离划分超平面越远则对该样本的分类越可靠，反之就不那么可靠。

那么，什么样的划分超平面是最佳超平面呢？

对于图1有A、B、C三个超平面，很明显应该选择超平面B，也就是说超平面首先应该能满足将两类样本点分开。

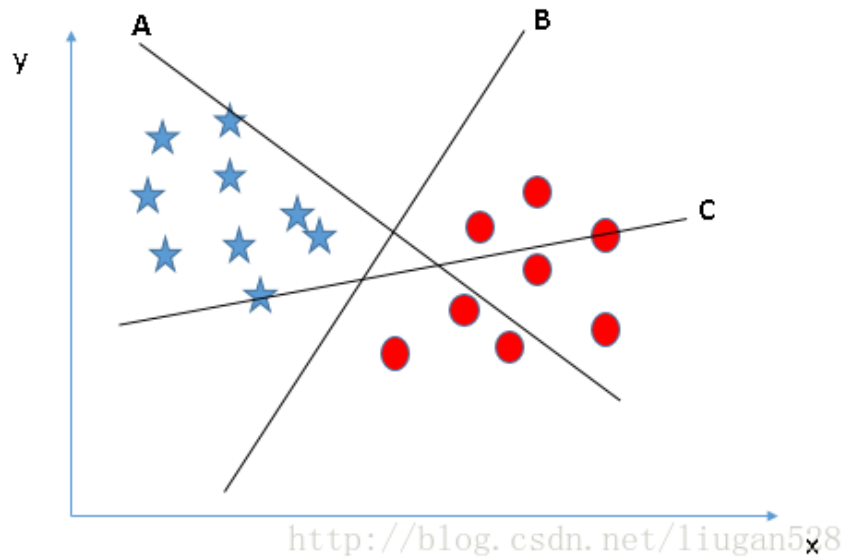


图1

对于图2的A、B、C三个超平面，应该选择超平面C，因为使用超平面C进行划分对训练样本局部扰动的“容忍”度最好，分类的鲁棒性最强。例如，由于训练集的局限性或噪声的干扰，训练集外的样本可能比图2中的训练样本更接近两个类目前的分隔界，在分类决策的时候就会出现错误，而超平面C受影响最小，也就是说超平面C所产生的分类结果是最鲁棒性的、是最可信的，对未见样本的泛化能力最强。

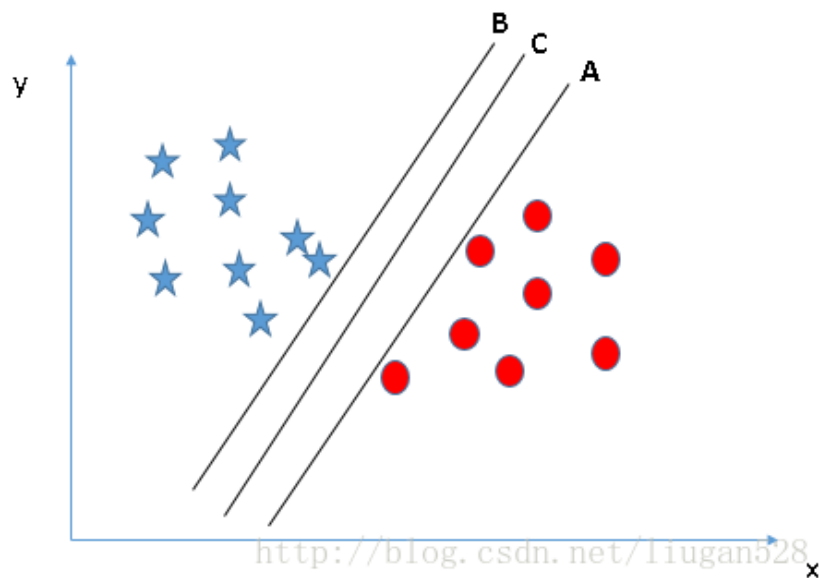


图2

下面以图3中示例进行推导得出最佳超平面。

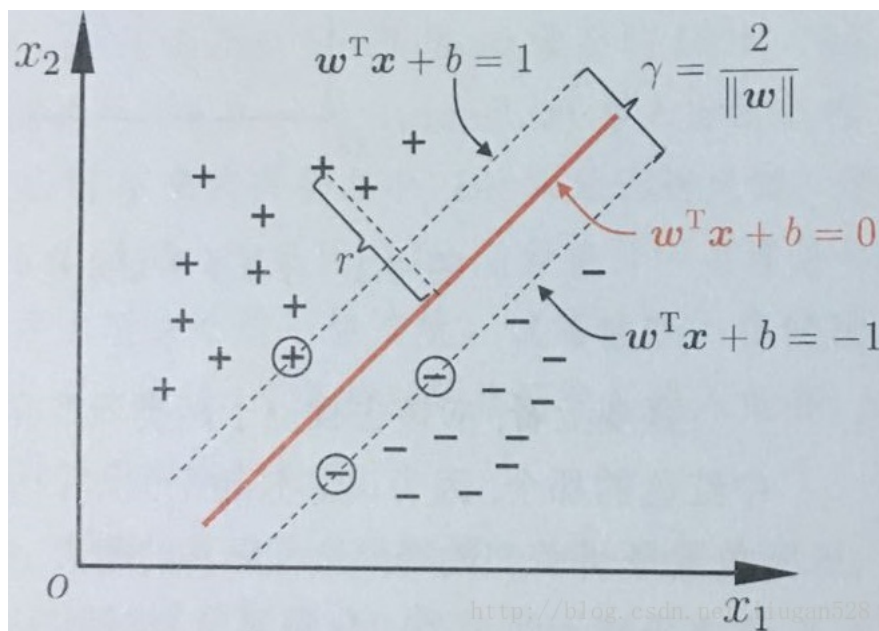


图3

空间中超平面可记为 $(\vec{w}, b)$  ( $w \rightarrow, b$ )，根据点到平面的距离公式，空间中任意点 $\vec{x}$ 到超平面 $(\vec{w}, b)$  ( $w \rightarrow, b$ )的距离可写为：

$$r = \frac{\vec{w}^T \vec{x} + b}{\|\vec{w}\|} \quad (3)$$

$$(3) r = \frac{\vec{w}^T \vec{x} + b}{\|\vec{w}\|}$$

假设超平面 $(\vec{w}, b)$  ( $w \rightarrow, b$ )能将训练样本正确分类，那么对于正样本一侧的任意一个样本

$(\vec{x}_i, y_i) \in D$  ( $x_i \rightarrow, y_i \in D$ )，应该需要满足该样本点往超平面的法向量 $\vec{w}$ 的投影到原点的距离大于一定值 $c$ 的时候使得该样本点被预测为正样本一类，即存在数值 $c$ 使得当 $\vec{w}^T \vec{x}_i > c$ 时

$y_i = +1$ 。  $\vec{w}^T \vec{x}_i > c$ 又可写为 $\vec{w}^T \vec{x}_i + b > 0$ 。在训练的时候我们要求限制条件更严格点以使最终得到的分类器鲁棒性更强，所以我们要求 $\vec{w}^T \vec{x}_i + b > 1$ 。也可以写为大于其它距离，但都可以通过同比例缩放 $\vec{w}$ 和 $b$ 来使得使其变为1，因此为计算方便这里直接选择1。同样对于负样本应该有 $\vec{w}^T \vec{x}_i + b < -1$ 。即：

$$\begin{cases} \vec{w}^T \vec{x}_i + b \geq +1, & y_i = +1 \\ \vec{w}^T \vec{x}_i + b \leq -1, & y_i = -1 \end{cases} \quad (4)$$

$$(4) \{ \vec{w}^T \vec{x}_i + b \geq +1, y_i = +1; \vec{w}^T \vec{x}_i + b \leq -1, y_i = -1 \}$$

亦即：

$$y_i (\vec{w}^T \vec{x}_i + b) \geq +1 \quad (5)$$

如图3所示，距离最佳超平面 $\vec{w}^T \vec{x} + b = 0$ 最近的几个训练样本点使上式中的等号成立，它们被称为“支持向量” (support vector)。记超平面 $\vec{w}^T \vec{x} + b = +1$ 和 $\vec{w}^T \vec{x} + b = -1$ 之间的距离为 $\gamma$ ，该距离又被称为“间隔” (margin)，SVM的核心之一就是想办法将“间隔” $\gamma$ 最大化。下面我们推导一下 $\gamma$ 与哪些因素有关：

记超平面 $\vec{w}^T \vec{x} + b = +1$ 上的正样本为 $\vec{x}_+$ ，超平面 $\vec{w}^T \vec{x} + b = -1$ 上的负样本为 $\vec{x}_-$ ，则根据向量的加减法规则 $\vec{x}_+$ 减去 $\vec{x}_-$ 得到的向量在最佳超平面的法向量 $\vec{w}$ 方向的投影即为“间隔” $\gamma$ ：

$$\gamma = (\vec{x}_+ - \vec{x}_-) \cdot \frac{\vec{w}}{\|\vec{w}\|} = \frac{\vec{x}_+ \cdot \vec{w}}{\|\vec{w}\|} - \frac{\vec{x}_- \cdot \vec{w}}{\|\vec{w}\|} \quad (6)$$

$$(6) \gamma = (\vec{x}_+ - \vec{x}_-) \cdot \frac{\vec{w}}{\|\vec{w}\|} = \frac{\vec{x}_+ \cdot \vec{w}}{\|\vec{w}\|} - \frac{\vec{x}_- \cdot \vec{w}}{\|\vec{w}\|}$$

而 $\vec{w}^T \vec{x}_+ + b = +1$ ， $\vec{w}^T \vec{x}_- + b = -1$ ，即：

$$\begin{cases} \vec{w}^T \vec{x}_+ = 1 - b \\ \vec{w}^T \vec{x}_- = -1 - b \end{cases} \quad (7)$$

$$(7) \{ \vec{w}^T \vec{x}_+ = 1 - b, \vec{w}^T \vec{x}_- = -1 - b$$

将(7)带入(6)可得：

$$\gamma = \frac{2}{\|\vec{w}\|} \quad (8)$$

$$(8) \gamma = 2 \|\vec{w}\|$$

也就是说使两类样本距离最大的因素仅仅和最佳超平面的法向量有关！

要找到具有“最大间隔”（maximum margin）的最佳超平面，就是找到能满足式(4)中约束的参数 $\vec{w}$ 和 $b$ 使得 $\gamma$ 最大，即：

$$\begin{cases} \max_{\vec{w}, b} \frac{2}{\|\vec{w}\|} \\ \text{s.t. } y_i (\vec{w}^T \vec{x}_i + b) \geq +1, i = 1, 2, \dots, n \end{cases} \quad (9)$$

$$(9) \{ \max_{\vec{w}, b} \frac{2}{\|\vec{w}\|} \text{s.t. } y_i (\vec{w}^T \vec{x}_i + b) \geq +1, i = 1, 2, \dots, n$$

显然(9)等价于

$$\begin{cases} \min_{\vec{w}, b} \frac{1}{2} \|\vec{w}\|^2 \\ \text{s.t. } y_i (\vec{w}^T \vec{x}_i + b) \geq +1, i = 1, 2, \dots, n \end{cases} \quad (10)$$

$$(10) \{ \min_{\vec{w}, b} \frac{1}{2} \|\vec{w}\|^2 \text{s.t. } y_i (\vec{w}^T \vec{x}_i + b) \geq +1, i = 1, 2, \dots, n$$

这就是SVM的基本型。

## 2.1 拉格朗日对偶问题

根据SVM的基本型求解出 $\vec{w}$ 和 $b$ 即可得到最佳超平面对应的模型：

$$f(\vec{x}) = \text{sign}(\vec{w}^T \vec{x} + b) \quad (11)$$

$$(11) f(\vec{x}) = \text{sign}(\vec{w}^T \vec{x} + b)$$

该求解问题本身是一个凸二次规划（convex quadratic programming）问题，可以通过开源的优化计算包进行求解，但是这样就无法体现SVM的精髓，我们可以将该凸二次规划问题通过拉格朗日对偶性来解决。对于式(10)的每条约束添加拉格朗日乘子 $\alpha_i \geq 0$ ，则该问题的拉格朗日函数可写为：

$$L(\vec{w}, b, \vec{\alpha}) = \frac{1}{2} \|\vec{w}\|^2 - \sum_{i=1}^n \alpha_i (y_i (\vec{w}^T \vec{x}_i + b) - 1) \quad (12)$$

$$(12) L(\vec{w}, b, \vec{\alpha}) = \frac{1}{2} \|\vec{w}\|^2 - \sum_{i=1}^n \alpha_i (y_i (\vec{w}^T \vec{x}_i + b) - 1)$$

其中 $\vec{\alpha} = (\alpha_1, \alpha_2, \dots, \alpha_n)$  分别是对应各个样本的拉格朗日乘子。

将 $L(\vec{w}, b, \vec{\alpha})$ 对 $\vec{w}$ 和 $b$ 求偏导并将偏导数等于零可得：

$$\begin{cases} \vec{w} = \sum_{i=1}^n \alpha_i y_i \vec{x}_i \\ \sum_{i=1}^n \alpha_i y_i = 0 \end{cases} \quad (13)$$

$$(13) \{ \vec{w} = \sum_{i=1}^n \alpha_i y_i \vec{x}_i, \sum_{i=1}^n \alpha_i y_i = 0$$

将(13)带入(12)消去 $\vec{w}$ 和 $b$ 就可得到(10)的对偶问题：

$$\begin{cases} \max_{\alpha} \sum_{i=1}^n \alpha_i - \frac{1}{2} \sum_{i=1}^n \sum_{j=1}^n \alpha_i \alpha_j y_i y_j \vec{x}_i^T \vec{x}_j \\ \text{s.t. } \alpha_i \geq 0, i = 1, 2, \dots, n \\ \sum_{i=1}^n \alpha_i y_i = 0 \end{cases} \quad (14)$$

$$(14) \{ \max_{\alpha} \sum_{i=1}^n \alpha_i - \frac{1}{2} \sum_{i=1}^n \sum_{j=1}^n \alpha_i \alpha_j y_i y_j \vec{x}_i^T \vec{x}_j \text{s.t. } \alpha_i \geq 0, i = 1, 2, \dots, n, \sum_{i=1}^n \alpha_i y_i = 0$$

由(14)可知我们并不关心单个样本是如何的，我们只关心样本间两两的乘积，这也为后面核方法提供了很大的便利。

求解出 $\vec{\alpha}$ 之后，再求解出 $\vec{w}$ 和 $b$ 即可得到SVM决策模型：

$$f(\vec{x}) = \vec{w}^T \vec{x} + b = \sum_{i=1}^n \alpha_i y_i \vec{x}_i^T \vec{x} + b \quad (15)$$

$$(15) f(x \rightarrow) = w \rightarrow T x \rightarrow + b = \sum_{i=1}^n \alpha_i y_i x_i \rightarrow T x \rightarrow + b$$

## 2.2 SVM问题的KKT条件

在(10)中有不等式约束，因此上述过程满足Karush-Kuhn-Tucker(KKT)条件：

$$\begin{cases} \alpha_i \geq 0 \\ y_i (w^T x + b) - 1 \geq 0 \\ \alpha_i (y_i (w^T x + b) - 1) = 0 \end{cases}, \quad i = 1, 2, \dots, n \quad (16)$$

$$(16) \{ \alpha_i \geq 0, y_i (w \rightarrow T x \rightarrow + b) - 1 \geq 0, i = 1, 2, \dots, n, \alpha_i (y_i (w \rightarrow T x \rightarrow + b) - 1) = 0$$

对于任意样本  $(\vec{x}_i, y_i)$  ( $x_i \rightarrow, y_i$ ) 总有  $\alpha_i = 0$  或  $y_i (w^T x + b) - 1 = 0$ 。如果  $\alpha_i = 0$  则  $y_i (w^T x + b) - 1 > 0$ ，表明对应的样本点在最大间隔边界上，即对应着支持向量。由此得出了SVM的一个重要性质：**训练完成之后，大部分的训练样本都不需要保留，最终的模型仅与支持向量有关。**

那么对于式(14)该如何求解  $\vec{\alpha}$  呢？很明显这是一个二次规划问题，可使用通用的二次规划算法来求解，但是SVM的算法复杂度是  $O(n^2)$ ，在实际问题中这种开销太大了。为了有效求解该二次规划问题，人们通过利用问题本身的特性，提出了很多高效算法，Sequential Minimal Optimization(SMO)就是一个常用的高效算法。在利用SMO算法进行求解的时候就需要用到上面的KKT条件。利用SMO算法求出  $\vec{\alpha}$  之后根据：

$$\begin{cases} \vec{w} = \sum_{i=1}^n \alpha_i y_i \vec{x}_i \\ y_i (w^T x + b) - 1 = 0 \end{cases} \quad (17)$$

$$(17) \{ w \rightarrow = \sum_{i=1}^n \alpha_i y_i x_i \rightarrow, y_i (w \rightarrow T x \rightarrow + b) - 1 = 0$$

即可求出  $\vec{w}$  和  $b$ 。求解出  $\vec{w}$  和  $b$  之后就可利用

$$f(x) = \text{sign}(w^T x + b) \quad (18)$$

$$(18) f(x \rightarrow) = \text{sign}(w \rightarrow T x \rightarrow + b)$$

进行预测分类了，注意在测试的时候不需要  $-1-1$ ，测试时没有训练的时候要求严格。

## 3. 软间隔支持向量机

在现实任务中很难找到一个超平面将不同类别的样本完全划分开，即很难找到合适的核函数使得训练样本在特征空间中线性可分。退一步说，即使找到了一个可以使训练集在特征空间中完全分开的核函数，也很难确定这个线性可分的结果是不是由于过拟合导致的。解决该问题的办法是在一定程度上运行SVM在一些样本上出错，为此引入了“软间隔” (soft margin) 的概念，如图4所示：

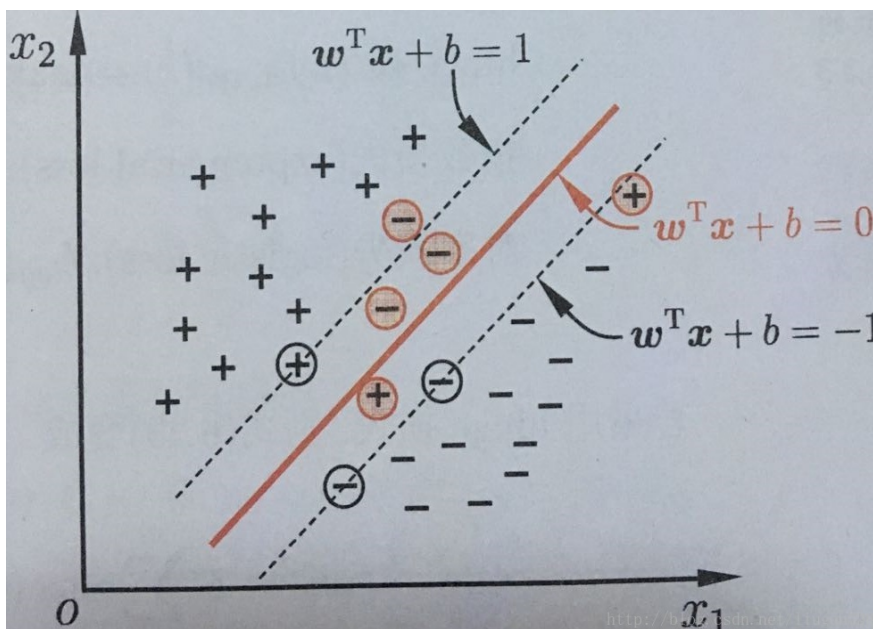


图4

具体来说，硬间隔支持向量机要求所有的样本均被最佳超平面正确划分，而软间隔支持向量机允许某些样本点不满足间隔大于等于1的条件 $y_i(\vec{w}\vec{x}_i + b) \geq 1$ ，当然在最大化间隔的时候也要限制不满足间隔大于等于1的样本的个数使之尽可能的少。于是我们引入一个惩罚系数 $C > 0$ ，并对每个样本点 $(\vec{x}_i, y_i)$ 引入一个松弛变量 (slack variables)  $\xi_i \geq 0$ ，此时可将式(10)改写为

$$\begin{cases} \min_{\vec{w}, b} (\frac{1}{2} \|\vec{w}\|^2 + C \sum_{i=1}^n \xi_i) \\ \text{s.t. } y_i(\vec{w}^T \vec{x}_i + b) \geq 1 - \xi_i, i = 1, 2, \dots, n \\ \xi_i \geq 0 \end{cases} \quad (19)$$

$$(19) \{ \min_{\vec{w}, b} (\frac{1}{2} \|\vec{w}\|^2 + C \sum_{i=1}^n \xi_i) \text{s.t. } y_i(\vec{w}^T \vec{x}_i + b) \geq 1 - \xi_i, i = 1, 2, \dots, n, \xi_i \geq 0$$

上式中约束条件改为 $y_i(\vec{w}\vec{x}_i + b) \geq 1 - \xi_i$ ，表示间隔加上松弛变量大于等于1；优化目标改为 $\min_{\vec{w}, b} (\frac{1}{2} \|\vec{w}\|^2 + C \sum_{i=1}^n \xi_i)$ ，表示对每个松弛变量都要有一个代价损失 $C\xi_i$ ， $C$ 越大对误分类的惩罚越大， $C$ 越小对误分类的惩罚越小。

式(19)是软间隔支持向量机的原始问题。可以证明 $\vec{w}$ 的解是唯一的， $b$ 的解不是唯一的， $b$ 的解是在一个区间内。假设求解软间隔支持向量机间隔最大化问题得到的最佳超平面是

$\vec{w}^* \vec{x} + b^* = 0$ ，对应的分类决策函数为

$$\begin{aligned} f(\vec{x}) &= \text{sign}(\vec{w}^* \vec{x} + b^*) \\ (20) f(\vec{x}) &= \text{sign}(\vec{w}^* \vec{x} + b^*) \end{aligned} \quad (20)$$

$f(\vec{x})$  称为软间隔支持向量机。

类似式(12)利用拉格朗日乘子法可得到上式的拉格朗日函数

$$L(\vec{w}, b, \vec{\alpha}, \vec{\xi}, \vec{\mu}) = \frac{1}{2} \|\vec{w}\|^2 + C \sum_{i=1}^n \xi_i - \sum_{i=1}^n \alpha_i (y_i(\vec{w}^T \vec{x}_i + b) - 1 + \xi_i) - \sum_{i=1}^n \mu_i \xi_i \quad (21)$$

$$(21) L(\vec{w}, b, \vec{\alpha}, \vec{\xi}, \vec{\mu}) = \frac{1}{2} \|\vec{w}\|^2 + C \sum_{i=1}^n \xi_i - \sum_{i=1}^n \alpha_i (y_i(\vec{w}^T \vec{x}_i + b) - 1 + \xi_i) - \sum_{i=1}^n \mu_i \xi_i$$

其中 $\alpha_i \geq 0$ ， $\mu_i \geq 0$ 是拉格朗日乘子。

令 $L(\vec{w}, b, \vec{\alpha}, \vec{\xi}, \vec{\mu})$ 分别对 $\vec{w}$ ， $b$ ， $\vec{\xi}$ 求偏导并将偏导数为零可得：

$$\begin{cases} \vec{w} = \sum_{i=1}^n \alpha_i y_i \vec{x}_i \\ \sum_{i=1}^n \alpha_i y_i = 0 \\ C = \alpha_i + \mu_i \end{cases} \quad (22)$$

$$(22) \{ \vec{w} = \sum_{i=1}^n \alpha_i y_i \vec{x}_i, \sum_{i=1}^n \alpha_i y_i = 0, C = \alpha_i + \mu_i$$

将式(22)带入式(21)便可得到式(19)的对偶问题：

$$\begin{cases} \max_{\vec{\alpha}} \sum_{i=1}^n \alpha_i - \frac{1}{2} \sum_{i=1}^n \sum_{j=1}^n \alpha_i \alpha_j y_i y_j \vec{x}_i^T \vec{x}_j \\ \text{s.t. } \sum_{i=1}^n \alpha_i y_i = 0, i = 1, 2, \dots, n \\ 0 \leq \alpha_i \leq C \end{cases} \quad (23)$$

$$(23) \{ \max_{\vec{\alpha}} \sum_{i=1}^n \alpha_i - \frac{1}{2} \sum_{i=1}^n \sum_{j=1}^n \alpha_i \alpha_j y_i y_j \vec{x}_i^T \vec{x}_j \text{s.t. } \sum_{i=1}^n \alpha_i y_i = 0, i = 1, 2, \dots, n, 0 \leq \alpha_i \leq C$$

对比软间隔支持向量机的对偶问题和硬间隔支持向量机的对偶问题可发现二者的唯一差别就在于对偶变量的约束不同，软间隔支持向量机对对偶变量的约束是 $0 \leq \alpha_i \leq C$ ，硬间隔支持向量机对对偶变量的约束是 $0 \leq \alpha_i$ ，于是可采用和硬间隔支持向量机相同的解法求解式(23)。同理在引入核方法之后同样能得到与式(23)同样的支持向量展开式。

类似式(16)对于软间隔支持向量机，KKT条件要求：

$$\begin{cases} \alpha_i \geq 0, \mu_i \geq 0 \\ y_i(\vec{w}^T \vec{x}_i + b) - 1 + \xi_i \geq 0 \\ \alpha_i (y_i(\vec{w}^T \vec{x}_i + b) - 1 + \xi_i) = 0 \\ \xi_i \geq 0, \mu_i \xi_i = 0 \end{cases} \quad (24)$$

$$(24) \{ \alpha_i \geq 0, \mu_i \geq 0, y_i(\vec{w}^T \vec{x}_i + b) - 1 + \xi_i \geq 0, \alpha_i (y_i(\vec{w}^T \vec{x}_i + b) - 1 + \xi_i) = 0, \xi_i \geq 0, \mu_i \xi_i = 0$$

同硬间隔支持向量机类似，对任意训练样本 $(\vec{x}_i, y_i)$ ，总有 $\alpha_i = 0$ 或

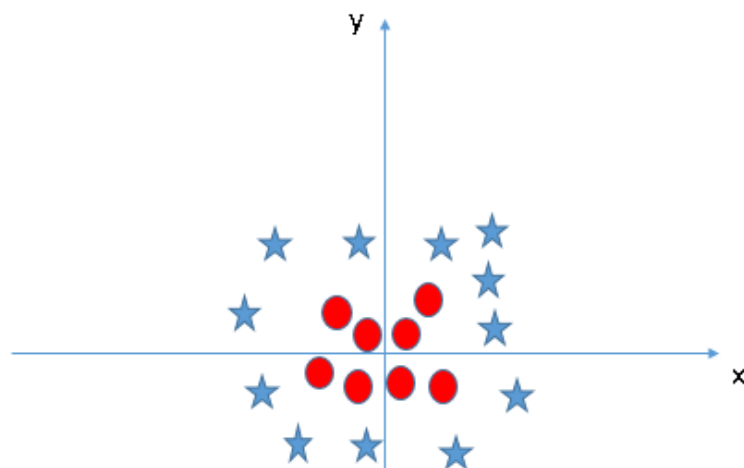
$y_i(\vec{w}^T \vec{x}_i + b) - 1 + \xi_i = 0$ ，若 $\alpha_i = 0$ ，则该样本不会对最佳决策面有任何影响；若

$\alpha_i > 0$  则必有  $y_i(\vec{w}\vec{x} + b) = 1 - \xi_i$ ，也就是说该样本是支持向量。由式(22)可知若  $\alpha_i < C$  则  $\mu_i > 0$  进而有  $\xi_i = 0$ ，即该样本处在最大间隔边界上；若  $\alpha_i = C$  则  $\mu_i = 0$  此时如果  $\xi_i \leq 1$  则该样本处于最大间隔内部，如果  $\xi_i > 1$  则该样本处于最大间隔外部即被分错了。由此也可看出，软间隔支持向量机的最终模型仅与支持向量有关。

## 4. 非线性支持向量机

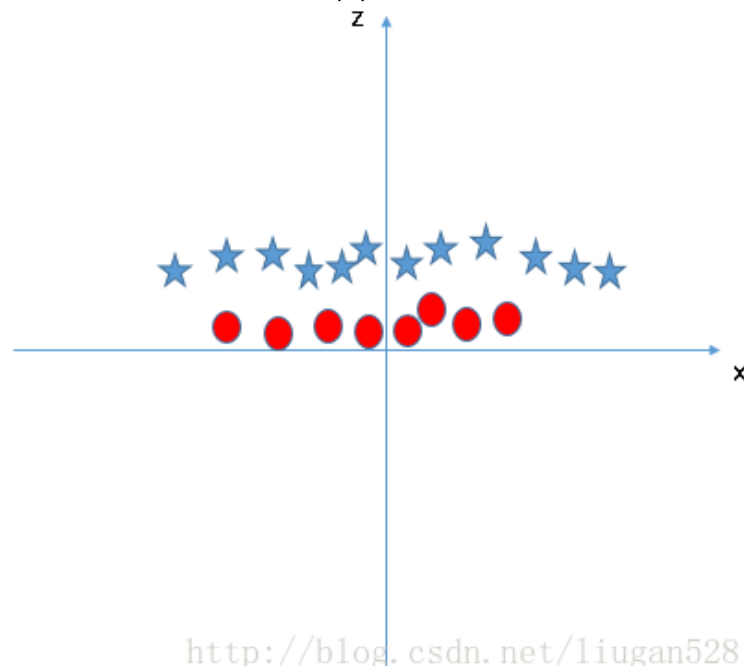
现实任务中原始的样本空间  $D$  中很可能并不存在一个能正确划分两类样本的超平面。例如图4中所示的问题就无法找到一个超平面将两类样本进行很好的划分。

对于这样的问题可以通过将样本从原始空间映射到特征空间使得样本在映射后的特征空间里线性可分。例如对图5做特征映射  $z = x^2 + y^2$  可得如图6所示的样本分布，这样就很好进行线性划分了。



<http://blog.csdn.net/liugan528>

图5



<http://blog.csdn.net/liugan528>

图6

令  $\phi(\vec{x})$  表示将样本点  $\vec{x}$  映射后的特征向量，类似于线性可分支持向量机中的表示方法，在特征空间



中划分超平面所对应的模型可表示为

$$f(\vec{x}) = \vec{w}^T \vec{x} + b \quad (25)$$

$$(25) f(\vec{x} \rightarrow) = \vec{w} \rightarrow^T \vec{x} + b$$

其中 $\vec{w}$ 和 $b$ 是待求解的模型参数。类似式(10)，有

$$\begin{cases} \min_{\vec{w}, b} \frac{1}{2} \|\vec{w}\|^2 \\ \text{s.t. } y_i (\vec{w}^T \vec{\phi}(\vec{x}) + b) \geq 1, i = 1, 2, \dots, n \end{cases} \quad (26)$$

$$(26) \{ \min_{\vec{w} \rightarrow, b} \frac{1}{2} \|\vec{w} \rightarrow\|^2 \text{ s.t. } y_i (\vec{w} \rightarrow^T \vec{\phi}(\vec{x} \rightarrow) + b) \geq 1, i = 1, 2, \dots, n \}$$

其拉格朗日对偶问题是

$$\begin{cases} \max_{\alpha} \sum_{i=1}^n \alpha_i - \frac{1}{2} \sum_{i=1}^n \sum_{j=1}^n \alpha_i \alpha_j y_i y_j \vec{\phi}(\vec{x}_i)^T \vec{\phi}(\vec{x}_j) \\ \text{s.t. } \alpha_i \geq 0, i = 1, 2, \dots, n \\ \sum_{i=1}^n \alpha_i y_i = 0 \end{cases} \quad (27)$$

$$(27) \{ \max_{\alpha} \sum_{i=1}^n \alpha_i - \frac{1}{2} \sum_{i=1}^n \sum_{j=1}^n \alpha_i \alpha_j y_i y_j \vec{\phi}(\vec{x}_i \rightarrow)^T \vec{\phi}(\vec{x}_j \rightarrow) \text{ s.t. } \alpha_i \geq 0, i = 1, 2, \dots, n, \sum_{i=1}^n \alpha_i y_i = 0 \}$$

求解(27)需要计算 $\vec{\phi}(\vec{x}_i)^T \vec{\phi}(\vec{x}_j)$ ，即样本映射到特征空间之后的内积，由于特征空间可能维度很高，甚至可能是无穷维，因此直接计算 $\vec{\phi}(\vec{x}_i)^T \vec{\phi}(\vec{x}_j)$ 通常是很困难的，在上文中我们提到其实我们根本不关心单个样本的表现，只关心特征空间中样本间两两的乘积，因此我们没有必要把原始空间的样本一个个地映射到特征空间中，只需要想办法求解出样本对应到特征空间中样本间两两的乘积即可。为了解决该问题可设想存在核函数：

$$\kappa(\vec{x}_i, \vec{x}_j) = \vec{\phi}(\vec{x}_i)^T \vec{\phi}(\vec{x}_j) \quad (28)$$

$$(28) \kappa(\vec{x}_i \rightarrow, \vec{x}_j \rightarrow) = \vec{\phi}(\vec{x}_i \rightarrow)^T \vec{\phi}(\vec{x}_j \rightarrow)$$

也就是说 $\vec{x}_i$ 与 $\vec{x}_j$ 在特征空间的内积等于它们在原始空间中通过函数 $\kappa(\cdot, \cdot)$ 计算的结果，这给求解带来很大的方便。于是式(27)可写为：

$$\begin{cases} \max_{\alpha} \sum_{i=1}^n \alpha_i - \frac{1}{2} \sum_{i=1}^n \sum_{j=1}^n \alpha_i \alpha_j y_i y_j \kappa(\vec{x}_i, \vec{x}_j) \\ \text{s.t. } \alpha_i \geq 0, i = 1, 2, \dots, n \\ \sum_{i=1}^n \alpha_i y_i = 0 \end{cases} \quad (29)$$

$$(29) \{ \max_{\alpha} \sum_{i=1}^n \alpha_i - \frac{1}{2} \sum_{i=1}^n \sum_{j=1}^n \alpha_i \alpha_j y_i y_j \kappa(\vec{x}_i \rightarrow, \vec{x}_j \rightarrow) \text{ s.t. } \alpha_i \geq 0, i = 1, 2, \dots, n, \sum_{i=1}^n \alpha_i y_i = 0 \}$$

同样的我们只关心在高维空间中样本之间两两乘积的结果而不关心样本是如何变换到高维空间中去的。求解后即可得到

$$f(\vec{x}) = \vec{w}^T \vec{\phi}(\vec{x}) + b = \sum_{i=1}^n \alpha_i y_i \vec{\phi}(\vec{x})^T \vec{\phi}(\vec{x}) + b = \sum_{i=1}^n \alpha_i y_i \kappa(\vec{x}, \vec{x}_i) + b \quad (30)$$

$$(30) f(\vec{x} \rightarrow) = \vec{w} \rightarrow^T \vec{\phi}(\vec{x} \rightarrow) + b = \sum_{i=1}^n \alpha_i y_i \vec{\phi}(\vec{x} \rightarrow)^T \vec{\phi}(\vec{x} \rightarrow) + b = \sum_{i=1}^n \alpha_i y_i \kappa(\vec{x} \rightarrow, \vec{x}_i \rightarrow) + b$$

剩余的问题同样是求解 $\alpha_i$ ，然后求解 $\vec{w}$ 和 $b$ 即可得到最佳超平面。

## 支持向量回归

支持向量机不仅可以用来解决分类问题还可以用来解决回归问题，称为支持向量回归（Support Vector Regression, SVR）。

对于样本 $(\vec{x}, y)$ 通常根据模型的输出 $f(\vec{x})$ 与真实值（即groundtruth） $y$ 之间的差别来计算损失，当且仅当 $f(\vec{x}) = y$ 时损失才为零。SVR的基本思路是允许预测值 $f(\vec{x})$ 与 $y$ 之间最多有 $\epsilon$ 的偏差，当 $|f(\vec{x}) - y| \leq \epsilon$ 时认为预测正确不计算损失，仅当 $|f(\vec{x}) - y| > \epsilon$ 时才计算损失。SVR问题可描述为：

$$\min_{\vec{w}, b} \left( \frac{1}{2} \|\vec{w}\|^2 + C \sum_{i=1}^n L_{\epsilon}(f(\vec{x}) - y_i) \right) \quad (31)$$

$$(31) \min_{\vec{w} \rightarrow, b} \left( \frac{1}{2} \|\vec{w} \rightarrow\|^2 + C \sum_{i=1}^n L_{\epsilon}(f(\vec{x} \rightarrow) - y_i) \right)$$

其中， $C \geq 0$ 为惩罚项， $L_{\epsilon}$ 为损失函数，定义为：

$$L_{\epsilon}(z) = \begin{cases} 0 & |z| \leq \epsilon \\ |z| - \epsilon & \text{otherwise} \end{cases} \quad (32)$$



$$(32) L\varepsilon(z) = \begin{cases} 0, & |z| \leq \varepsilon \\ |z| - \varepsilon, & \text{otherwise} \end{cases}$$

进一步地引入松弛变量  $\xi_i, \hat{\xi}_i$ , 则新的最优化问题为:

$$\begin{cases} \min_{\vec{w}, b, \xi, \hat{\xi}} \left( \frac{1}{2} \|\vec{w}\|^2 + C \sum_{i=1}^n (\xi_i + \hat{\xi}_i) \right) \\ \text{s.t.} \quad \begin{aligned} f(\vec{x}_i) - y_i &\leq \varepsilon + \xi_i \\ y_i - f(\vec{x}_i) &\leq \varepsilon + \hat{\xi}_i \\ \xi_i &\geq 0, \hat{\xi}_i \geq 0 \end{aligned} \end{cases}, i = 1, 2, \dots, n \quad (33)$$

$$(33) \begin{cases} \min_{\vec{w}, b, \xi, \hat{\xi}} \left( \frac{1}{2} \|\vec{w}\|^2 + C \sum_{i=1}^n (\xi_i + \hat{\xi}_i) \right) \\ \text{s.t.} \quad \begin{aligned} f(\vec{x}_i) - y_i &\leq \varepsilon + \xi_i \\ y_i - f(\vec{x}_i) &\leq \varepsilon + \hat{\xi}_i \\ \xi_i &\geq 0, \hat{\xi}_i \geq 0 \end{aligned} \end{cases}, i = 1, 2, \dots, n$$

这就是SVR的原始问题。类似地引入拉格朗日乘子  $\mu_i \geq 0, \hat{\mu}_i \geq 0, \alpha_i \geq 0, \hat{\alpha}_i \geq 0$ , 则对应的拉格朗日函数为:

$$L(\vec{w}, b, \vec{\alpha}, \vec{\hat{\alpha}}, \vec{\xi}, \vec{\hat{\xi}}, \vec{\mu}, \vec{\hat{\mu}}) = \frac{1}{2} \|\vec{w}\|^2 + C \sum_{i=1}^n (\xi_i + \hat{\xi}_i) - \sum_{i=1}^n \mu_i \xi_i - \sum_{i=1}^n \hat{\mu}_i \hat{\xi}_i + \sum_{i=1}^n \alpha_i (f(\vec{x}_i) - y_i - \varepsilon - \xi_i) + \sum_{i=1}^n \hat{\alpha}_i (y_i - f(\vec{x}_i) - \varepsilon - \hat{\xi}_i) \quad (34)$$

$$(34) \begin{cases} L(\vec{w}, b, \vec{\alpha}, \vec{\hat{\alpha}}, \vec{\xi}, \vec{\hat{\xi}}, \vec{\mu}, \vec{\hat{\mu}}) = \frac{1}{2} \|\vec{w}\|^2 + C \sum_{i=1}^n (\xi_i + \hat{\xi}_i) \\ - \sum_{i=1}^n \mu_i \xi_i - \sum_{i=1}^n \hat{\mu}_i \hat{\xi}_i + \sum_{i=1}^n \alpha_i (f(\vec{x}_i) - y_i - \varepsilon - \xi_i) + \sum_{i=1}^n \hat{\alpha}_i (y_i - f(\vec{x}_i) - \varepsilon - \hat{\xi}_i) \end{cases}$$

令  $L(\vec{w}, b, \vec{\alpha}, \vec{\hat{\alpha}}, \vec{\xi}, \vec{\hat{\xi}}, \vec{\mu}, \vec{\hat{\mu}})$  对  $\vec{w}, b, \vec{\xi}, \vec{\hat{\xi}}, \vec{\mu}, \vec{\hat{\mu}}$  的偏导数为零可得:

$$\begin{cases} \vec{w} = \sum_{i=1}^n (\alpha_i - \hat{\alpha}_i) \vec{x}_i \\ \sum_{i=1}^n (\alpha_i - \hat{\alpha}_i) = 0 \\ C = \alpha_i + \mu_i \\ C = \hat{\alpha}_i + \hat{\mu}_i \end{cases} \quad (35)$$

$$(35) \begin{cases} \vec{w} = \sum_{i=1}^n (\alpha_i - \hat{\alpha}_i) \vec{x}_i \\ \sum_{i=1}^n (\alpha_i - \hat{\alpha}_i) = 0 \\ C = \alpha_i + \mu_i \\ C = \hat{\alpha}_i + \hat{\mu}_i \end{cases}$$

将式(35)代入式(34)即可得到SVR的对偶问题:

$$\begin{cases} \max_{\alpha, \hat{\alpha}} \sum_{i=1}^n (y_i (\alpha_i - \hat{\alpha}_i) - \varepsilon (\alpha_i + \hat{\alpha}_i) - \frac{1}{2} \sum_{i=1}^n \sum_{j=1}^n (\alpha_i - \hat{\alpha}_i) (\alpha_j - \hat{\alpha}_j) \vec{x}_i^T \vec{x}_j) \\ \text{s.t.} \quad \sum_{i=1}^n (\alpha_i - \hat{\alpha}_i) = 0 \quad 0 \leq \alpha_i, \hat{\alpha}_i \leq C \end{cases} \quad (36)$$

$$(36) \begin{cases} \max_{\alpha, \hat{\alpha}} \sum_{i=1}^n (y_i (\alpha_i - \hat{\alpha}_i) - \varepsilon (\alpha_i + \hat{\alpha}_i) - \frac{1}{2} \sum_{i=1}^n \sum_{j=1}^n (\alpha_i - \hat{\alpha}_i) (\alpha_j - \hat{\alpha}_j) \vec{x}_i^T \vec{x}_j) \\ \text{s.t.} \quad \sum_{i=1}^n (\alpha_i - \hat{\alpha}_i) = 0, 0 \leq \alpha_i, \hat{\alpha}_i \leq C \end{cases}$$

其KKT条件为:

$$\begin{cases} \alpha_i (f(\vec{x}_i) - y_i - \varepsilon - \xi_i) = 0 \\ \hat{\alpha}_i (y_i - f(\vec{x}_i) - \varepsilon - \hat{\xi}_i) = 0 \\ \alpha_i \hat{\alpha}_i = 0, \xi_i \hat{\xi}_i = 0 \\ (C - \alpha_i) \xi_i = 0, (C - \hat{\alpha}_i) \hat{\xi}_i = 0 \end{cases} \quad (37)$$

$$(37) \begin{cases} \alpha_i (f(\vec{x}_i) - y_i - \varepsilon - \xi_i) = 0 \\ \hat{\alpha}_i (y_i - f(\vec{x}_i) - \varepsilon - \hat{\xi}_i) = 0 \\ \alpha_i \hat{\alpha}_i = 0, \xi_i \hat{\xi}_i = 0 \\ (C - \alpha_i) \xi_i = 0, (C - \hat{\alpha}_i) \hat{\xi}_i = 0 \end{cases}$$

SVR的解形如:

$$f(\vec{x}) = \sum_{i=1}^n (\alpha_i - \hat{\alpha}_i) \vec{x}_i^T \vec{x} + b \quad (38)$$

$$(38) f(\vec{x}) = \sum_{i=1}^n (\alpha_i - \hat{\alpha}_i) \vec{x}_i^T \vec{x} + b$$

进一步地如果引入核函数则SVR可表示为:

$$f(\vec{x}) = \sum_{i=1}^n (\alpha_i - \hat{\alpha}_i) \kappa(\vec{x}_i, \vec{x}) + b \quad (39)$$

$$(39) f(\vec{x}) = \sum_{i=1}^n (\alpha_i - \hat{\alpha}_i) \kappa(\vec{x}_i, \vec{x}) + b$$

其中  $\kappa(\vec{x}_i, \vec{x}) = \phi(\vec{x}_i)^T \phi(\vec{x})$   $\kappa(\vec{x}_i, \vec{x}) = \phi(\vec{x}_i)^T \phi(\vec{x})$  为核函数。

## 常用核函数

名称	表达式	参数
线性核	$\kappa(\vec{x}_i, \vec{x}_j) = \vec{x}_i^T \vec{x}_j$ $\kappa(x_i \rightarrow, x_j \rightarrow) = x_i \rightarrow^T x_j \rightarrow$	
多项式核	$\kappa(\vec{x}_i, \vec{x}_j) = (\vec{x}_i^T \vec{x}_j)^n$ $\kappa(x_i \rightarrow, x_j \rightarrow) = (x_i \rightarrow^T x_j \rightarrow)^n$	$n \geq 1$ $n \geq 1$ 为多项式的次数
高斯核 (RBF)	$\kappa(\vec{x}_i, \vec{x}_j) = \exp\left(-\frac{\ \vec{x}_i - \vec{x}_j\ ^2}{2\sigma^2}\right)$ $\kappa(x_i \rightarrow, x_j \rightarrow) = \exp(-\ x_i \rightarrow - x_j \rightarrow\ ^2 / 2\sigma^2)$	$\sigma > 0$ $\sigma > 0$ 为高斯核的带宽
拉普拉斯核	$\kappa(\vec{x}_i, \vec{x}_j) = \exp\left(-\frac{\ \vec{x}_i - \vec{x}_j\ }{\sigma}\right)$ $\kappa(x_i \rightarrow, x_j \rightarrow) = \exp(-\ x_i \rightarrow - x_j \rightarrow\  / \sigma)$	$\sigma > 0$
Sigmoid核	$\kappa(\vec{x}_i, \vec{x}_j) = \tanh(\beta \vec{x}_i^T \vec{x}_j + \theta)$ $\kappa(x_i \rightarrow, x_j \rightarrow) = \tanh(\beta x_i \rightarrow^T x_j \rightarrow + \theta)$	$\tanh$ 为双曲正切函数

## 5. SVM的优缺点

优点：

SVM在中小量样本规模的时候容易得到数据和特征之间的非线性关系，可以避免使用神经网络结构选择和局部极小值问题，可解释性强，可以解决高维问题。

缺点：

SVM对缺失数据敏感，对非线性问题没有通用的解决方案，核函数的正确选择不容易，计算复杂度高，主流的算法可以达到 $O(n^2)$ 的复杂度，这对大规模的数据是吃不消的。

## 6. 参考文献

周志华. 机器学习 [D]. 清华大学出版社, 2016.

华校专、王正林. Python大战机器学习 [D]. 电子工业出版社, 2017.

Peter Flach著、段菲译. 机器学习 [D]. 人民邮电出版社, 2016.

[Understanding Support Vector Machine algorithm from examples \(along with code\).](#)

[KKT条件介绍](#)



“分享知识，共同进步！”

GaryLau 的赞赏码

更多资料请移步github：

<https://github.com/GarryLau/MachineLearning>