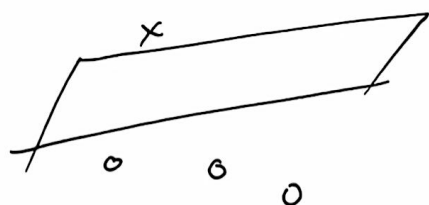


SVM 推导与系统:

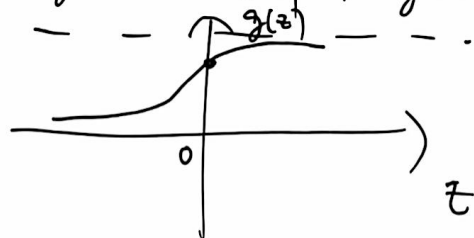
SVM support vector machine.

在 classifier 中, 我们会确定一条 boundary, $h_\theta(x)$ 会与这个 boundary 密切相关. 在高维空间中, 这个 boundary 就是一个 hyperplane.



因此, 我们可以将这个面作为一个重要的决策依据。

在 sigmoid 函数中, $g(z)$ 形如



z 由 $\theta^T x$ 拟合,

直观感受是, $|\theta^T x|$ 越大, 即离 $z=0$ 越远, 我们越有把握认为 $\hat{y} = g(\theta^T x)$ 处于其中某一种状态, 假如 $|\theta^T x|$ 很小, 我们判断越模糊。

因此引入新的优化目标:

(我们此时 ~~将 y 限制在 $\{0, 1\}$~~ 而是 $y \in \{-1, 1\}$)

我们希望 ~~$y=0$ 时~~ $y=-1$ 时

~~y 尽可能~~ $\theta^T x$ 尽可能小 (为负)

$y=1$ 时 $\theta^T x$ 尽可能大 (为正)

因此提出 functional margin

$$\gamma = y(\theta^T x)$$

如果将 $x_0=1$ 项分离出来,

$$\gamma = y(w^T x + b)$$

但这会引出一个问题, 我们知道

当 $\lambda \neq 0$ 时 $\lambda(w^T x + b) = 0$ 是同一个超平面, 但 $\gamma = y(w^T x + b)$

$$= \lambda y(w^T x + b)$$

我们可以仅通过不断放缩 (w^T, b)

提高 γ , 这能实现优化的“目标” (提高 γ) 但达不到目的 (提高预测的精度)

我们可以通过规范化 w^T ~~来~~ 克服上述的漏洞。

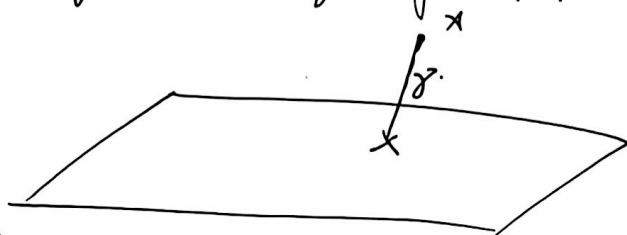
$$\hat{\gamma} = y \left(\frac{w^T}{\|w\|} x + \frac{b}{\|w\|} \right)$$

此时并缩放 (w^T, b) 将不再影响 $\hat{\gamma}$

同时, 此时 $\frac{w^T x + b}{\|w\|}$ 就刚好是

geometric margin.

geometric margin 计算.



首先, ~~w^T~~ \vec{w} 总是与 $w^T x + b = 0$ 正交的。

假设 $w^T x + b$ 上有两点 x_1, x_2

$$\begin{cases} w^T x_1 + b = 0 \\ w^T x_2 + b = 0 \end{cases} \Rightarrow w^T (x_1 - x_2) = 0$$

这说明 \vec{w} 与 $w^T x + b = 0$ 正交。



因此, 假设 x 在 $w^T x + b$ 上投影的点是 p .

$$\vec{x} - \vec{p} = \lambda \frac{\vec{w}}{\|\vec{w}\|}$$

$$\vec{p} = \vec{x} - \lambda \frac{\vec{w}}{\|\vec{w}\|} \quad \leftarrow \text{一个单位向量}$$

$$\vec{p} \text{ 满足 } w^T x_p + b = 0.$$

$$w^T (x - \lambda \frac{\vec{w}}{\|\vec{w}\|}) + b = 0.$$

$$w^T x - \lambda \frac{w^T w}{\|\vec{w}\|} + b = 0.$$

$$\lambda = \frac{w^T x + b}{\|\vec{w}\|}$$

然后, 我们有 SVM 的核心算法.

optimal margin classifier.

我们在考虑一组样本 $(x^{(i)}, y^{(i)}) \sim D$ 时, 我们希望找到一个 hyperplane, 将 $\{x^{(i)} | y^{(i)} = 1\}$ 和 $\{x^{(i)} | y^{(i)} = -1\}$ 分离开。此时, 我们要选取一个合适的 γ , 这样我们才能更好地迭代优化。

在 SVM 中, 我们不再关心样本数据所服从的某种分布, 而是关心算法如何才能在

train set 上有更好的预测的把握。

更甚, 我们更加关心那些离 hyperplane

更近的点, 我们更希望在这种边缘上上 S.t. $y^{(i)} (w^T x^{(i)} + b) \geq 1$ 有更好的准确率。

于是, 我们想优化下面的问题:

$$\gamma = \min \gamma$$

我们要通过调整 w 与 b $\max \gamma$,

$$\text{即 } \max_{w, b, \gamma} \gamma$$

$$\text{s.t. } \forall_i \gamma \leq \gamma$$

$$\text{即 } \max_{w, b} \gamma$$

$$\text{s.t. } \forall_i \frac{y^{(i)} (w^T x^{(i)} + b)}{\|\vec{w}\|} \geq \gamma$$

$$\text{即 } \max_{w, b, \gamma} \gamma$$

$$\text{s.t. } \forall_i y^{(i)} (w^T x^{(i)} + b) \geq \gamma$$

$$\|\vec{w}\| = 1$$

$\|\vec{w}\| = 1$ 的约束条件并不好实现, 我们希望约束条件是一个 convex 的函数或者 affine 的平面, 这样可以很好地使用

Lagrange dual.

$$\text{即 } \max_{w, b, \gamma} \left(\frac{\gamma}{\|\vec{w}\|} \right)$$

geometric margin.

$$\forall_i y^{(i)} (w^T x^{(i)} + b) \geq \gamma$$

functional margin.

因为 γ 是与 $\|\vec{w}\|$ 密切相关的, 我们总是可以通过 scale $\|\vec{w}\|$ 找到合适的 (w, b) 满足 $\gamma = 1$, 且这不影响 $\frac{\gamma}{\|\vec{w}\|}$ (geometric margin)

$$\max_{w, b, \gamma} \frac{1}{\|\vec{w}\|}$$

$$\Leftrightarrow \min_{w, b, \gamma} \frac{1}{2} \|\vec{w}\|_2^2$$

$$\text{s.t. } \forall_i y^{(i)} (w^T x^{(i)} + b) \geq 1$$

注意, 此时, 我们已经达到目标:

优化的函数 $\frac{1}{2} \|\vec{w}\|_2^2$ 是 convex 的, 约束条件是 affine 的



于是, 我们有

$$\mathcal{L}(w, b, \alpha) = \frac{1}{2} \|w\|_2^2 + \sum_{i=1}^m \alpha_i (-y^{(i)} (w^T x^{(i)} + b) + 1)$$

where $\forall i, \alpha_i \geq 0$.

primal:

$$\min_{w, b} \max_{\alpha} \mathcal{L}(w, b, \alpha)$$

dual:

$$\max_{\alpha} \min_{w, b} \mathcal{L}(w, b, \alpha)$$

我们先解决 dual 问题是 (后面再证明这个问题是符合 KKT 条件的 $\Leftrightarrow p^* = d^*$)

$$\begin{aligned} \nabla_w \mathcal{L}(w, b, \alpha) &= w + \sum_{i=1}^m \alpha_i (-y^{(i)} x^{(i)}) \\ &= w - \sum_{i=1}^m \alpha_i y^{(i)} x^{(i)} = 0 \quad (1) \end{aligned}$$

$$\begin{aligned} \nabla_b \mathcal{L}(w, b, \alpha) &= \sum_{i=1}^m \alpha_i (-y^{(i)}) = 0 \quad (2) \end{aligned}$$

$$\text{由 (1) 得 } w^* = \sum_{i=1}^m \alpha_i y^{(i)} x^{(i)} \quad (3)$$

将 (2) (3) 代入 $\mathcal{L}(w, b, \alpha)$

$$\mathcal{L}(w, b, \alpha)$$

$$\begin{aligned} \min_{w, b} \mathcal{L}(w, b, \alpha) &= \frac{1}{2} \sum_{i=1}^m \alpha_i y^{(i)} x^{(i)T} \sum_{j=1}^m \alpha_j y^{(j)} x^{(j)} \\ &+ \sum_{i=1}^m \alpha_i (-y^{(i)}) \sum_{j=1}^m \alpha_j y^{(j)} x^{(j)T} x^{(i)} + \sum_{i=1}^m \alpha_i \end{aligned}$$

这里有点混了, 重写一下

$$= \frac{1}{2} \sum_{i=1}^m \alpha_i y^{(i)} x^{(i)T} \sum_{j=1}^m \alpha_j y^{(j)} x^{(j)} + \sum_{i=1}^m \alpha_i (-y^{(i)}) \sum_{j=1}^m \alpha_j y^{(j)} x^{(j)T} x^{(i)} + \sum_{i=1}^m \alpha_i$$

整理下

$$\begin{aligned} \min_{w, b} \mathcal{L}(w, b, \alpha) &= \sum_{i=1}^m \alpha_i - \frac{1}{2} \sum_{i=1}^m \sum_{j=1}^m y^{(i)} y^{(j)} \alpha_i \alpha_j \langle x^{(i)}, x^{(j)} \rangle \end{aligned}$$

剩下的问题是

$$\max_{\alpha} \sum_{i=1}^m \alpha_i - \frac{1}{2} \sum_{i=1}^m \sum_{j=1}^m y^{(i)} y^{(j)} \alpha_i \alpha_j \langle x^{(i)}, x^{(j)} \rangle$$

$$\text{s.t. } \forall i, \alpha_i \geq 0, \sum_{i=1}^m \alpha_i y^{(i)} = 0$$

以上的算法是对那些离 boundary 很近的点上非常敏感的 (那些点就是 support vector)

这会带来一些问题

① 在拉 Lagrange 对偶问题中, $\max_{\alpha} \min_{w, b} \mathcal{L}(w, b, \alpha)$ 在限制被违反时, 会得到无穷

② 该算法受 support vector 的影响非常大, 一个表现有异常的点, 会使边界产生较大的变动。

所以会产生一定的想法, 容忍一些点有少量的违反条件:

$$\min_{w, b, \xi, \epsilon} \frac{1}{2} \|w\|_2^2 + C \sum \epsilon_i$$

$$\text{s.t. } y^{(i)} (w^T x^{(i)} + b) \geq 1 - \epsilon_i$$

for $i=1, 2, \dots, m$

$$\text{且 } \epsilon_i \geq 0 \text{ for } i=1, 2, \dots, m$$

这里的 C 是一个 hyperparameter, 做为惩罚。



得到 Lagrangian.

④

$$\mathcal{L}(w, b, \varepsilon, \alpha, \beta)$$

$$= \frac{1}{2} \|w\|_2^2 + C \sum_{i=1}^m \varepsilon_i + \sum_{i=1}^m \alpha_i (-y^{(i)} (w^T x^{(i)} + b) + 1 - \varepsilon_i) + \sum_{i=1}^m \beta_i (-\varepsilon_i)$$

$$\nabla_w \mathcal{L} = w + \sum_{i=1}^m \alpha_i (-y^{(i)}) x^{(i)} = 0 \quad ①$$

$$\nabla_b \mathcal{L} = \sum_{i=1}^m \alpha_i (-y^{(i)}) = 0 \quad ②$$

$$\frac{\partial}{\partial \varepsilon_i} \mathcal{L} = C + \alpha_i (-1) + \beta_i (-1) = 0. \quad ③$$

将 ① ② ③ 代入

$$W(\alpha) = \left(\sum_{i=1}^m \alpha_i - \frac{1}{2} \sum_{i,j=1}^m y^{(i)} y^{(j)} \alpha_i \alpha_j \langle x^{(i)}, x^{(j)} \rangle \right)$$

转化后的问题是

$$\max_{\alpha} \sum_{i=1}^m \alpha_i - \frac{1}{2} \sum_{i,j=1}^m y^{(i)} y^{(j)} \alpha_i \alpha_j \langle x^{(i)}, x^{(j)} \rangle$$

s.t.

$$0 \leq \alpha_i \leq C$$

$$\sum_{i=1}^m \alpha_i y^{(i)} = 0.$$

