

第2章 感知机模型

适用条件

解决什么问题

三要素

并
针对二分类问题. 假设为线性可分的

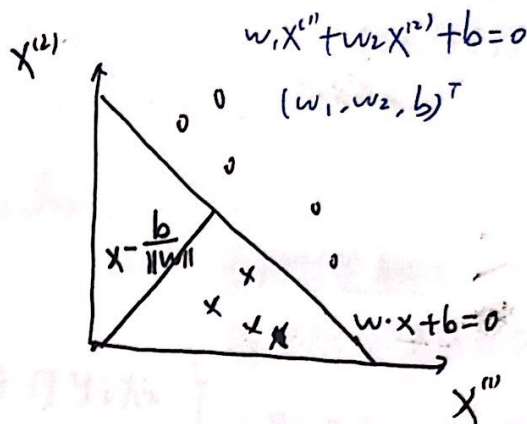
2.1 符号表示

输入空间 $X \subseteq R^n$

输入变量 $x \in X$

输出空间 $Y = \{+1, -1\}$

输出变量 $y \in \{+1, -1\}$



假设空间 \rightarrow 符号函数

$$f(x) = \text{sign}(w \cdot x + b)$$

$$= \begin{cases} +1, & w \cdot x + b \geq 0 \\ -1, & w \cdot x + b < 0 \end{cases}$$

$\therefore x$ 是个 n 维向量

$\therefore w$ 也是 n 维向量

$w \cdot x$ 是内积的意思

$$= w_1 x^{(1)} + w_2 x^{(2)} + \dots + w_n x^{(n)}$$

2.2 感知机学习策略

★

求 x_i 到一条直线的距离

损失函数: 输入 n 维
则用 $n-1$ 维超平面去分

误分类点到超平面的总距离

$$\frac{|w \cdot x_i + b|}{\|w\|}$$

二范式

$$\hookrightarrow \|w\| = \sqrt{w_1^2 + \dots + w_n^2}$$

$$L(w, b) = - \sum_{x_i \in M} y_i (w \cdot x_i + b)$$

推导
来原:

解出 w 与 b

$$\sum_{x_i \in M} \frac{|w \cdot x_i + b|}{\|w\|}$$

$$= \sum_{x_i \in M} - \frac{y_i (w \cdot x_i + b)}{\|w\|}$$

要使其最小, $\therefore \|w\|$ 为正数

$- \dots \pm b$

2.3 感知机学习算法 (目的: 求出 w, b , 使其为以下损失函数极小化的问题)

算法 2.1 (随机梯度下降法) $\min_{w, b} L(w, b) = -\sum_{x_i \in M} y_i (w \cdot x + b)$

输入: 训练数据集 $T = [(x_1, y_1), \dots, (x_n, y_n)]$, 学习率 η

1. 选取初值 w_0, b_0 ($0 < \eta \leq 1$)

2. 在训练集中选取数据 (x_i, y_i)

3. 如果 $y_i (w \cdot x_i + b) \leq 0$ 要同时更新
 $w := w_{\text{前}} + \eta y_i x_i$
 $b := b + \eta y_i$ } 每次都更靠近正确的
 (取决于初始的点与更新顺序)

4. 转至 2, 直到训练集中没有误分类的点

输出: w, b

感知机模型对偶形式

算法 2.2

$$f(x) = \text{sign} \left(\underbrace{\sum_{j=1}^N \alpha_j y_j x_j}_w x + b \right)$$

输入: 训练数据集 T ,

学习率 η

$\alpha = (\alpha_1, \dots, \alpha_N)^T$
 由 (α_j, b) 决定的 j 由 $1 \sim N$

1. 初值 $\alpha_i = 0$ $b_i = 0$

2. 在训练集中选取数据 (x_i, y_i)

3. 如果 $y_i \left(\sum_{j=1}^N \alpha_j y_j x_j \cdot x + b \right) \leq 0$

$$\alpha_i := \alpha_i + \eta$$

$$b := b + \eta y_i$$

4. 转至 2, 直至训练集中没有误分类点

★ 对偶算法与 ~~随机~~ 梯度下降算法 比较

① 对偶算法每次更新的时候只用更新 α_i 与 b 两个数
而

梯度下降算法每次都要更新一整个 w 向量与 b 这个数

② 判断是否分类错误时

对偶算法 算的是 x_i 与 w 的内积, 可以预先算好

↓ 计算量小

↓

即 Gram 矩阵

而

梯度下降算法需要计算一整个 w 向量

证:

定理 2.1 (Novikoff) 算法的收敛性

① 前提: 记 $\hat{w} = (w^T, b)^T$, $\hat{w} \in R^{n+1}$ $\hat{w} \cdot \hat{x} = w \cdot x + b$
 $\hat{x} = (x^T, 1)^T$, $\hat{x} \in R^{n+1}$

② \hat{w}_{opt} 意义解释

设训练数据集线性可分, 则超平面表示为 $w \cdot x + b = 0$

找到的这个超平面为 $w_{opt} \cdot x + b_{opt} = 0$

$$\begin{aligned} \text{则 } \hat{w}_{opt} = (w_{opt}^T, b_{opt})^T \\ \hat{x} = (x^T, 1)^T \end{aligned} \Rightarrow \hat{w}_{opt} \cdot \hat{x} = w_{opt} \cdot x + b_{opt}$$

\parallel
 0

$$\therefore \hat{w}_{opt} \cdot \hat{x} = 0$$

$$\therefore \boxed{\geq \hat{w}_{opt}} \hat{x} = 0$$

抽象为 \hat{w}
新的

为了让 \hat{w}_{opt} 比较唯一

我们约束 $\|\hat{w}_{opt}\| = 1$

念伽马

\therefore ① 我们要证明 存在 $\|\hat{w}_{opt}\| = 1$ 的超平面 $\hat{w}_{opt} \cdot \hat{x} = 0$

使得训练集完全分开, 且存在 $\gamma > 0$, 使得

$$y_i (\hat{w}_{opt} \cdot \hat{x}_i) \geq \gamma$$

证 ①

$\therefore w_{opt}$ 的平面使所有数据分开

\therefore 有 $y_i (\hat{w}_{opt} \cdot \hat{x}_i) > 0$ (对所有实例来说)

\therefore 让 $\gamma = \min_i y_i (\hat{w}_{opt} \cdot \hat{x}_i)$ 即 $\gamma \leq y_i (\hat{w}_{opt} \cdot \hat{x}_i)$

② 令 $R = \max_{1 \leq i \leq N} \|\hat{x}_i\|$, $k \leq \left(\frac{R}{\gamma}\right)^2$ 可以理解为这个向量 \hat{x}_i 长度的表示.
修正的次数 上一个定理中的 γ

分2步

1) 要证 $\hat{w}_k \cdot \hat{w}_{opt} \geq k\eta\gamma$ 念 yī tāi 念 qā mǎ

$$= (\hat{w}_{k-1} + \eta y_i \hat{x}_i) \cdot \hat{w}_{opt}$$

由上一个定理有 $y_i (\hat{w}_{opt} \cdot \hat{x}_i) \geq \gamma$

$$= \hat{w}_{k-1} \cdot \hat{w}_{opt} + \eta y_i \hat{x}_i \cdot \hat{w}_{opt}$$

$$\begin{aligned} \therefore & \geq \hat{w}_{k-1} \cdot \hat{w}_{opt} + \eta \gamma \\ & \geq \hat{w}_{k-2} \cdot \hat{w}_{opt} + \eta \gamma + \eta \gamma \\ & \vdots \\ & \geq \hat{w}_0 \cdot \hat{w}_{opt} + k\eta \gamma \end{aligned}$$

\therefore 假设 $\hat{w}_0 = (0, \dots, 0)^T$

$$\therefore = k\eta \gamma$$

a 与 b 都是向量

$$\|a+b\|^2 = (a+b)^T (a+b)$$

$$= \|a\|^2 + 2a \cdot b + \|b\|^2$$

2) 要证 $\|\hat{w}_k\|^2$ 要小于某个值

$$\begin{aligned} &= \|\hat{w}_{k-1} + \eta y_i \hat{x}_i\|^2 \\ &= \|\hat{w}_{k-1}\|^2 + 2\eta y_i \hat{w}_{k-1} \cdot \hat{x}_i + \eta^2 \|\hat{x}_i\|^2 \end{aligned}$$

$$\therefore \leq \|\hat{w}_{k-1}\|^2 + \eta^2 R^2$$

$$\leq \|\hat{w}_0\|^2 + k\eta^2 R^2$$

$$\therefore = k\eta^2 R^2$$

★ **重要**

对 i 来说这是个

误分类点

$$\therefore 2\eta y_i \hat{w}_{k-1} \cdot \hat{x}_i < 0$$

$$\therefore \text{假设 } \|\hat{x}_i\|^2 \leq R^2$$

$\therefore \hat{w}_0$ 为 $(0, \dots, 0)^T$

得出

∴ 由上面2个证明可得出

$$\begin{cases} \hat{w}_k \cdot \hat{w}_{opt} \geq k\eta\gamma & \text{--- ①} \\ \|\hat{w}_k\|^2 \leq k\eta^2 R^2 & \text{--- ②} \end{cases}$$

★ 素雅 $k \leq (\frac{R}{\gamma})^2$

★ 根据柯西不等式

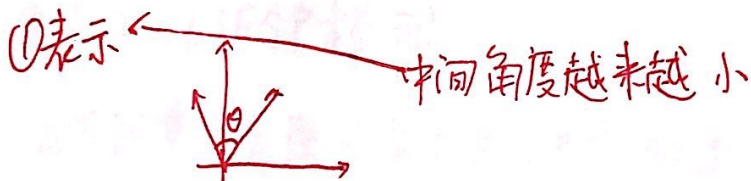
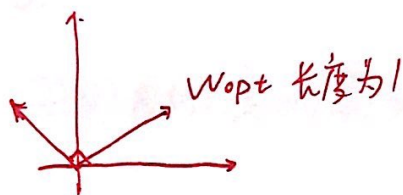
$$k\eta\gamma \leq \hat{w}_k \cdot \hat{w}_{opt} \leq \|\hat{w}_k\| \|\hat{w}_{opt}\| \stackrel{=1}{\leq} \sqrt{k\eta^2 R^2}$$

$$\text{即 } k\eta\gamma \leq \sqrt{k\eta^2 R^2} \Rightarrow k \leq (\frac{R}{\gamma})^2$$

★ 理解 ①, ② 两式

① 表示 $\hat{w}_k \cdot \hat{w}_{opt}$ 在逐渐增大, ∴ k 在 ↑

② 表示 $\|\hat{w}_k\|$ 长度可以被限制住



★ 即 ①② 联合起来大概表示为 w_k 在逐渐靠近 w_{opt}
(重要)