



机器学习

Machine Learning

第四讲：分类的基本学习算法



注：

本章讲义符号与前章略有不同，
用 t 表示标注， y 表示模型输出

本讲义采用两套符号之一：

用 y 表示标注， \hat{y} 表示模型输出
或：

用 t 表示标注， y 表示模型输出



1. 分类问题 (Classification)

- 数据集 (标注集)

$$\mathcal{D} = \{(\mathbf{x}_i, t_i)\}_{i=1}^N \quad \Rightarrow \mathbf{y} = h(\mathbf{x})$$

- 分类(Classification) $t, y \in \{1, \dots, C\}$

$C=2$, *binary classification*

$C > 2$, *multiclass classification*



基本分类问题表示

对于2类问题

$$t \in \{0, 1\}$$

$$t = 1 \Rightarrow C1 \quad t = 0 \Rightarrow C2$$

对于K (>2)类问题 (K-to-1编码)

$$t = (0, \dots, 0, 1, 0, \dots, 0)^T$$



分类的三种基本模型

(1). 判决函数模型

线性模型

$$y(\mathbf{x}) = \mathbf{w}^T \mathbf{x} + w_0$$

广义线性模型 (generalized linear models)

$$y(\mathbf{x}) = f\left(\mathbf{w}^T \mathbf{x} + w_0\right)$$

$f(\)$ 为激活函数

由训练集学习参数 \mathbf{w}, w_0 (或用验证集确定超参数)

X 表示数据集



分类的三种基本模型（续）

(2). 判决（概率）模型

由数据集直接训练后验概率，由决策论确定输出

$$p(C_k | \mathbf{x}; \mathbf{X}) \quad \text{简写为} \quad p(C_k | \mathbf{x})$$



分类的三种基本模型（续）

(3). 生成（概率）模型

首先得到 $p(\mathbf{x}|C_k)$ 和 $p(C_k)$ 或: $p(\mathbf{x}, C_k)$

再由Bayes公式得

$$p(C_k|\mathbf{x}) = \frac{p(\mathbf{x}|C_k)p(C_k)}{p(\mathbf{x})}$$

$p(\mathbf{x})$ $p(\mathbf{x}, C_k)$ 可以求得，能够生成更多数据



2. 判决函数方法 (Discriminant Functions)

- 线性分类
 - (2类、多类) , LS优化
- Fisher线性判决函数
 - (2类、多类) (统计方法的传统技术)
 - 投影到低维空间, 投影可分辨力最大化
- 感知器算法 (The Perceptron Algorithm)
 - (MLP的最简化形式, 曾起到重要作用)
- 传统算法 (略, 详见教材4.2节)

3.逻辑回归 (Logistic Regression)

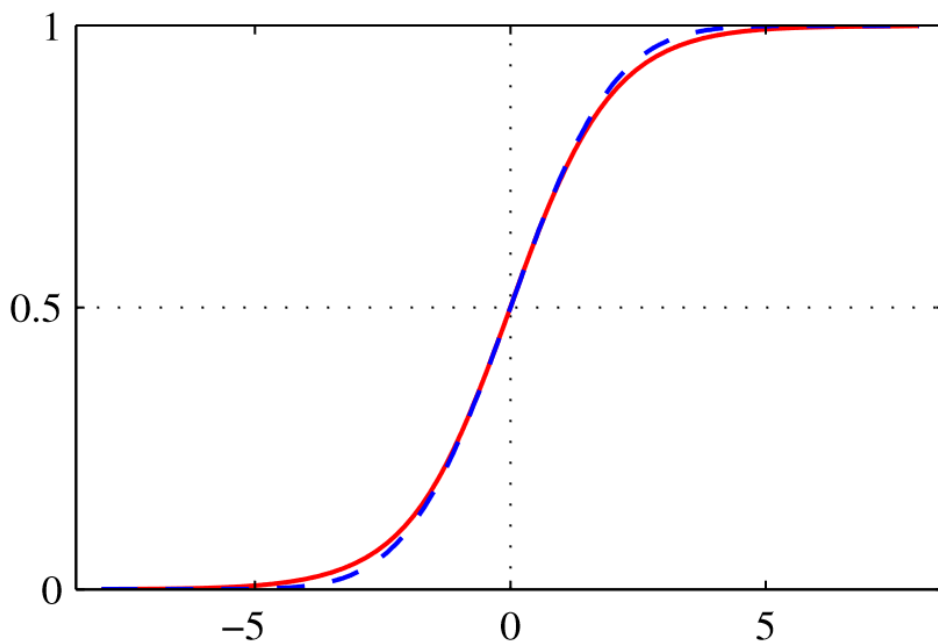
(概率判决模型)

3.1 二类问题

定义: logistics sigmoid函数

$$\sigma(a) = \frac{1}{1 + \exp(-a)}$$

用logistic sigmoid表示
类的后验概率



本节直接采用基函数 $\varphi(x)$
线性形式时取 $\varphi(x) = \bar{x}$



logistics sigmoid函数的性质

$$\sigma(-a) = 1 - \sigma(a)$$

$$\frac{d\sigma(a)}{da} = \sigma(a)(1 - \sigma(a))$$

对权向量的线性系统

$$a = \mathbf{w}^T \phi(\mathbf{x})$$

$$\mathbf{w} = [w_0, w_1, \dots, w_{M-1}]^T$$

其中

$$\phi(\mathbf{x}) = [1, \phi_1(\mathbf{x}), \phi_2(\mathbf{x}), \dots, \phi_{M-1}(\mathbf{x})]^T$$



二分类后验概率的表示

$$p(\mathcal{C}_1|\phi) = y(\phi) = \sigma(\mathbf{w}^T \phi)$$

另：

$$p(\mathcal{C}_2|\phi) = 1 - p(\mathcal{C}_1|\phi)$$

由训练样本直接学习参数 \mathbf{w}

对于新的输入 $\mathbf{x} \Rightarrow \phi(\mathbf{x})$

计算类后验概率 $P(\mathcal{C}_1|\phi) = y(\phi) = \sigma(\mathbf{w}^T \phi(\mathbf{x}))$ 并分类



二分类逻辑回归参数学习

训练样本集 $\{\mathbf{x}_n, t_n\}_{n=1}^N$

变换的样本集 $\{\phi(\mathbf{x}_n), t_n\}_{n=1}^N = \{\phi_n, t_n\}_{n=1}^N$

几个简写

$$\phi_n = \phi(\mathbf{x}_n)$$

$$\mathbf{t} = (t_1, \dots, t_N)^T$$

$$y_n = p(\mathcal{C}_1 | \phi_n) \quad y_n = \sigma(a_n)$$

$$a_n = \mathbf{w}^T \phi_n$$



二分类逻辑回归参数学习（续）

似然函数

$$p(\mathbf{t}|\mathbf{w}) = \prod_{n=1}^N y_n^{t_n} \{1 - y_n\}^{1-t_n}$$

负对数似然函数

$$E(\mathbf{w}) = -\ln p(\mathbf{t}|\mathbf{w})$$

$$= -\sum_{n=1}^N \{t_n \ln y_n + (1 - t_n) \ln(1 - y_n)\}$$

注意

$$y_n = \sigma(a_n)$$

含参数 \mathbf{w}

二分类逻辑回归参数学习（续）



目标函数对 \mathbf{w} 的梯度，可导出为

$$\begin{aligned}\nabla E(\mathbf{w}) &= \sum_{n=1}^N (y_n - t_n) \phi_n \\ &= \sum_{n=1}^N \left(\sigma(\mathbf{w}^T \phi(\mathbf{x}_n)) - t_n \right) \phi(\mathbf{x}_n)\end{aligned}$$

随机梯度

$$\nabla E_n = (y_n - t_n) \phi_n$$

二分类逻辑回归参数学习（续）

随机梯度算法学习参数

$$\begin{aligned}\mathbf{w}^{(k+1)} &= \mathbf{w}^{(k)} - \eta \nabla E_n \\ &= \mathbf{w}^{(k)} - \eta (y_n - t_n) \phi_n \\ &= \mathbf{w}^{(k)} - \eta \left(\sigma \left(\mathbf{w}^{(k)T} \phi(\mathbf{x}_n) \right) - t_n \right) \phi(\mathbf{x}_n)\end{aligned}$$

η 学习率参数

注：可按一定次序使用 $\{\mathbf{x}_n, t_n\}_{n=1}^N$ ，甚至可循环使用直到收敛，也可用小批量平均梯度。



二分类逻辑回归参数学习（续）

IRLS算法*

(Iterative Reweighted Least Squares)

利用牛顿迭代思想

$$\mathbf{w}^{(\text{new})} = \mathbf{w}^{(\text{old})} - \mathbf{H}^{-1} \nabla E(\mathbf{w})$$

其中

$$\mathbf{H} = \nabla \nabla E(\mathbf{w}) \quad \text{Hessian matrix}$$



IRLS算法

Φ $N \times M$ 矩阵

ϕ_n^T 第 n 行

$$\nabla E(\mathbf{w}) = \sum_{n=1}^N (y_n - t_n) \phi_n = \Phi^T (\mathbf{y} - \mathbf{t})$$

$$\begin{aligned} \mathbf{H} &= \nabla \nabla E(\mathbf{w}) \\ &= \sum_{n=1}^N y_n (1 - y_n) \phi_n \phi_n^T = \Phi^T \mathbf{R} \Phi \end{aligned}$$

其中 $N \times N$ diagonal matrix \mathbf{R} $R_{nn} = y_n(1 - y_n)$

\mathbf{H} 正定，有唯一最优解



IRLS算法

$$\begin{aligned}\mathbf{w}^{(\text{new})} &= \\ &= \mathbf{w}^{(\text{old})} - (\Phi^T \mathbf{R} \Phi)^{-1} \Phi^T (\mathbf{y} - \mathbf{t}) \\ &= (\Phi^T \mathbf{R} \Phi)^{-1} \{ \Phi^T \mathbf{R} \Phi \mathbf{w}^{(\text{old})} - \Phi^T (\mathbf{y} - \mathbf{t}) \} \\ &= (\Phi^T \mathbf{R} \Phi)^{-1} \Phi^T \mathbf{R} \mathbf{z}\end{aligned}$$

这里 $\mathbf{z} = \Phi \mathbf{w}^{(\text{old})} - \mathbf{R}^{-1}(\mathbf{y} - \mathbf{t})$ N维向量

这是一个 \mathbf{R} 作为加权矩阵的加权LS，每次需重新计算 \mathbf{R} ，重新运行LS，直到收敛。故名：IRLS。

抗逻辑回归的overfitting!

正则化逻辑回归

(Regularized Logistic Regression)

$$\begin{aligned} E(\mathbf{w}) &= -\ln p(\mathbf{t}|\mathbf{w}) + \frac{\lambda}{2} \mathbf{w}^T \mathbf{w} \\ &= -\sum_{n=1}^N \{t_n \ln y_n + (1 - t_n) \ln(1 - y_n)\} \\ &\quad + \frac{\lambda}{2} \mathbf{w}^T \mathbf{w} \end{aligned}$$

也可加更一般的正则化项，例如：
q=1（或<1）对应稀疏逻辑回归。

$$\frac{\lambda}{2} \sum_{j=1}^M |w_j|^q$$



3.2 多类逻辑回归

Multiclass logistic regression

定义Softmax 函数

$$p(\mathcal{C}_k|\phi) = y_k(\phi) = \frac{\exp(a_k)}{\sum_j \exp(a_j)}$$

这里 $a_k = \mathbf{w}_k^T \phi$

对每一类定义类后验概率

对每一类定义和学习权向量 \mathbf{w}_k



多类逻辑回归

Multiclass logistic regression

Softmax 函数的性质

$$\frac{\partial y_k}{\partial a_j} = y_k (I_{kj} - y_j)$$

这里

$$I_{jk} = \begin{cases} 1, & j = k \\ 0, & j \neq k \end{cases}$$



多类逻辑回归参数学习

训练样本集 $\{\mathbf{x}_n, \mathbf{t}_n\}_{n=1}^N$

这里

$$\mathbf{t}_n = [0, \dots, 0, 1, 0, \dots, 0]^T$$
$$= [t_{n1}, t_{n2}, \dots, t_{nK}]^T$$

变换的样本集 $\{\phi(\mathbf{x}_n), \mathbf{t}_n\}_{n=1}^N = \{\phi_n, \mathbf{t}_n\}_{n=1}^N$

令：

$$\mathbf{T} = [t_{nk}]_{N \times K} \quad y_{nk} = y_k(\phi_n)$$



多类逻辑回归参数学习

似然函数

$$\begin{aligned} p(\mathbf{T} | \mathbf{w}_1, \dots, \mathbf{w}_K) &= \\ &= \prod_{n=1}^N \prod_{k=1}^K p(\mathcal{C}_k | \phi_n)^{t_{nk}} = \prod_{n=1}^N \prod_{k=1}^K y_{nk}^{t_{nk}} \end{aligned}$$



多类逻辑回归参数学习

$$\begin{aligned} E(\mathbf{w}_1, \dots, \mathbf{w}_K) &= -\ln p(\mathbf{T} | \mathbf{w}_1, \dots, \mathbf{w}_K) \\ &= -\sum_{n=1}^N \sum_{k=1}^K t_{nk} \ln y_{nk} \end{aligned}$$

存在约束条件

$$\sum_k t_{nk} = 1$$



多类逻辑回归参数学习

目标函数对各参数向量的梯度

$$\nabla_{\mathbf{w}_j} E(\mathbf{w}_1, \dots, \mathbf{w}_K) = \sum_{n=1}^N (y_{nj} - t_{nj}) \phi_n$$

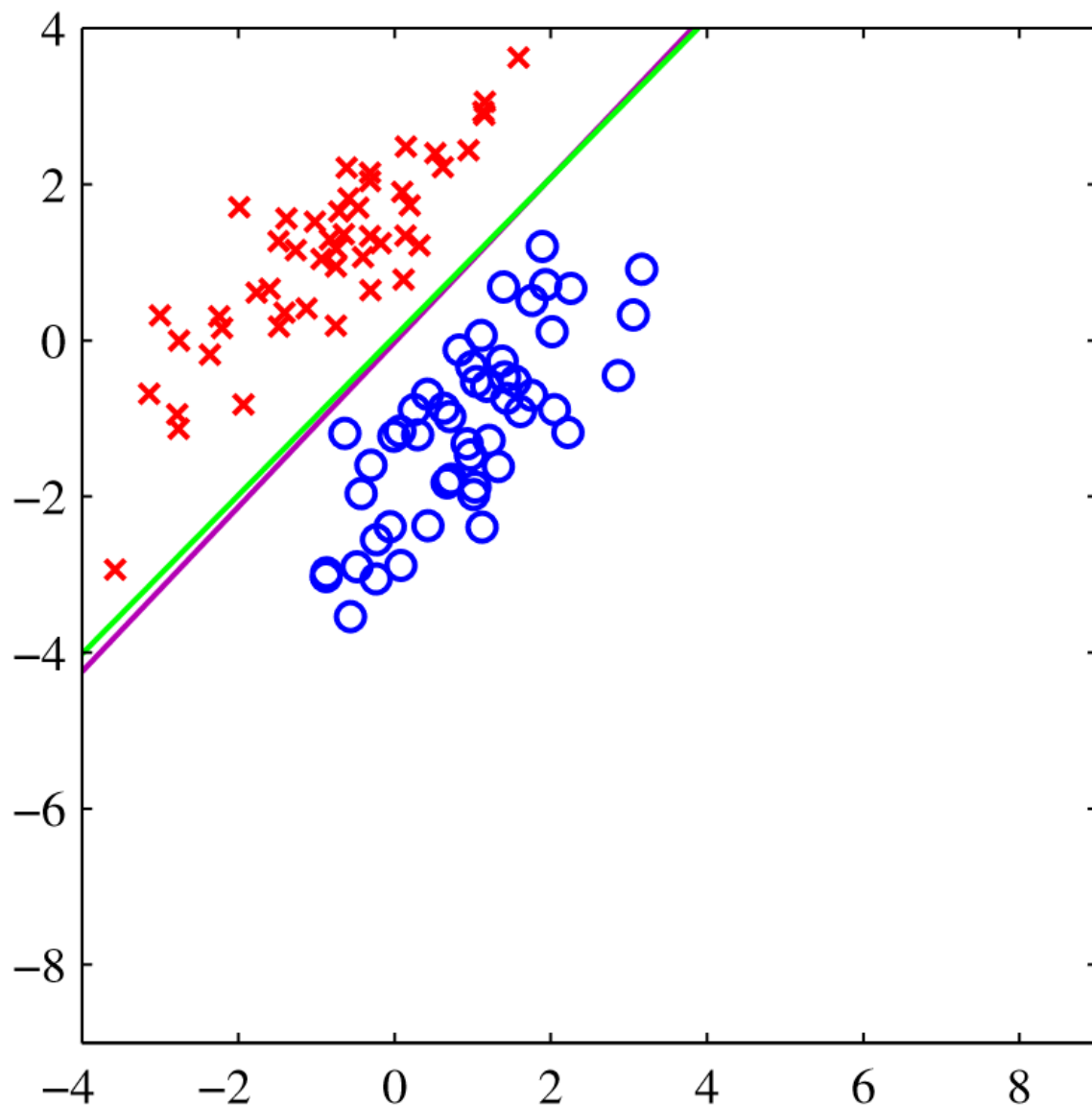
随机梯度

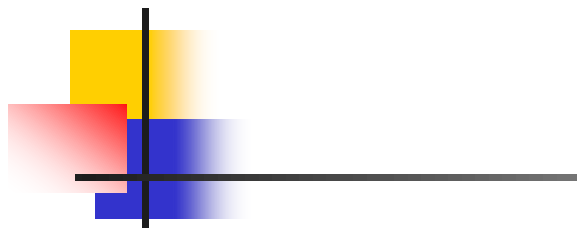
$$\nabla_{\mathbf{w}_j} E_n = (y_{nj} - t_{nj}) \phi_n$$

对每一个参数向量 \mathbf{w}_j ，可分别应用随机梯度算法迭代

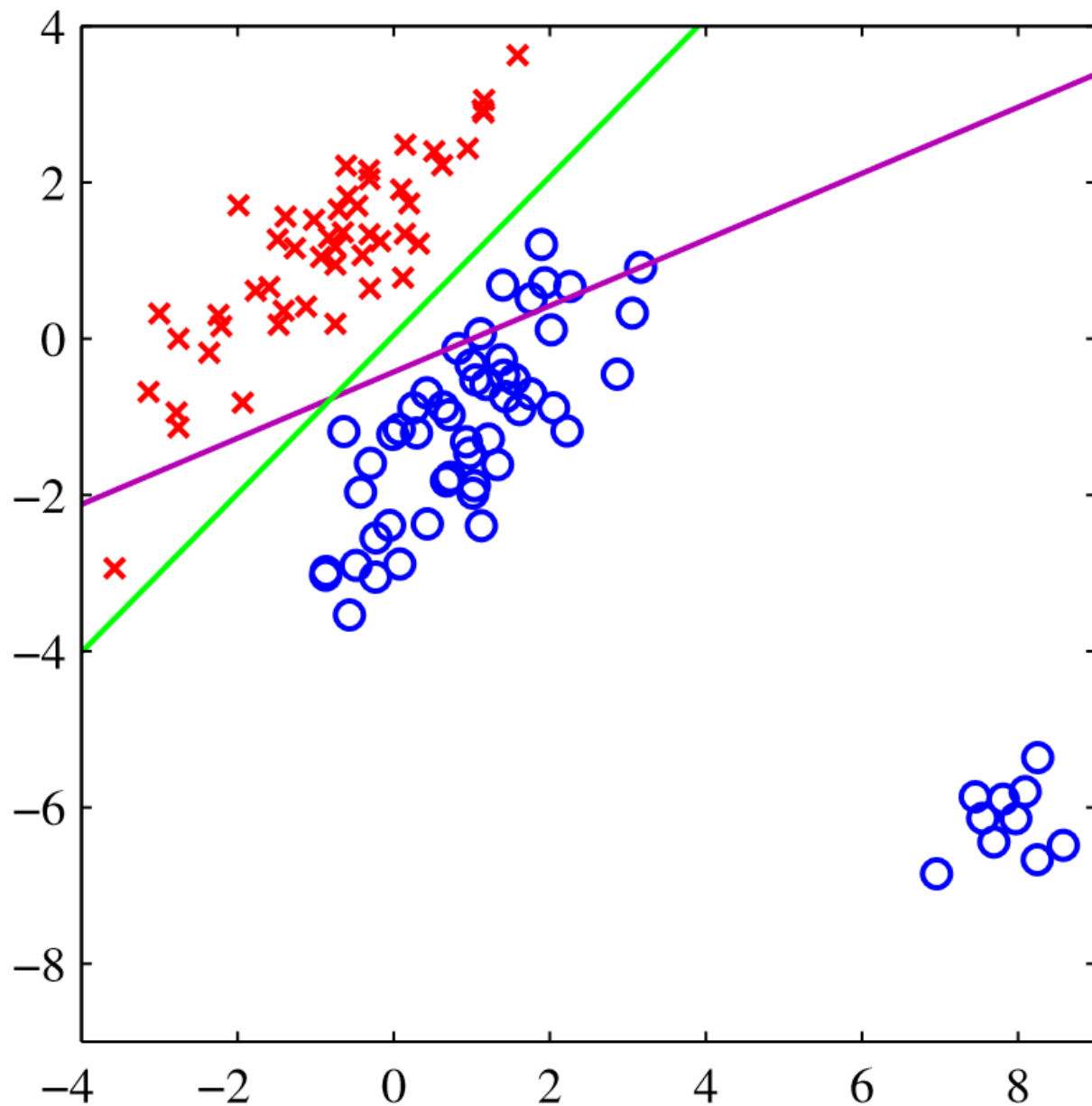
3.3 例

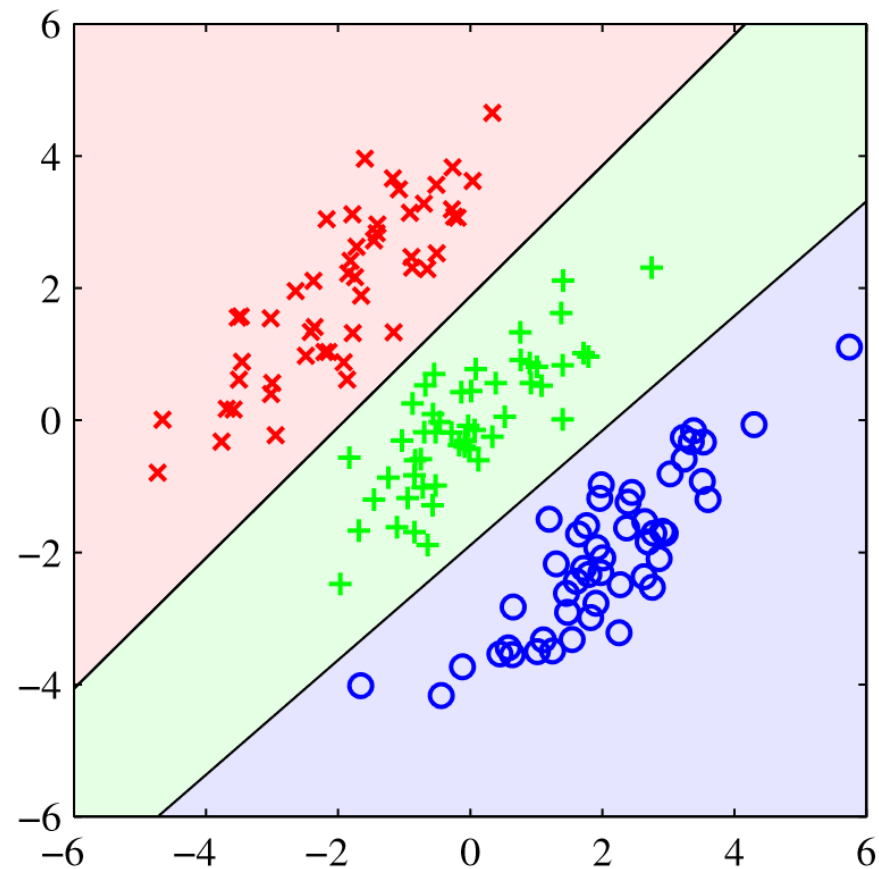
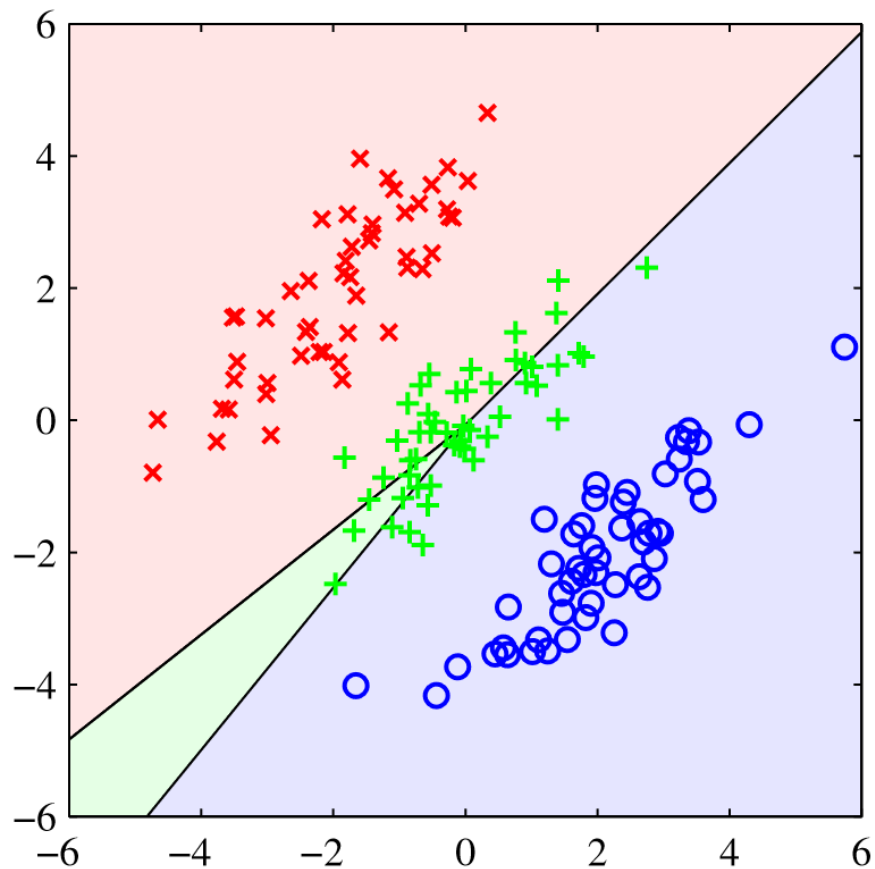
例子：两类，
分别用LS分类
和逻辑回归
绿线是逻辑回归
紫线是LS的边界





例子：右下角增加
几个野点，LS分类
明显变差，
逻辑回归基本不变





三类情况，即使是清晰可分的，LS分类也很差（右）
此例中，逻辑回归分类比较理想



4. 分类的生成模型

若可通过数据集模型化: $p(\mathbf{x}|\mathcal{C}_k)$ 、 $p(\mathcal{C}_k)$ 或 $p(x, t)$

类 \mathcal{C}_1 的后验概率可写为

$$\begin{aligned} p(\mathcal{C}_1|\mathbf{x}) &= \frac{p(\mathbf{x}|\mathcal{C}_1)p(\mathcal{C}_1)}{p(\mathbf{x}|\mathcal{C}_1)p(\mathcal{C}_1) + p(\mathbf{x}|\mathcal{C}_2)p(\mathcal{C}_2)} \\ &= \frac{1}{1 + \exp(-a)} = \sigma(a) \quad \# (1) \end{aligned}$$

$$\text{其中} \quad a = \ln \frac{p(\mathbf{x}|\mathcal{C}_1)p(\mathcal{C}_1)}{p(\mathbf{x}|\mathcal{C}_2)p(\mathcal{C}_2)} \quad \# (2)$$



4.1 离散生成模型: Naïve Bayes

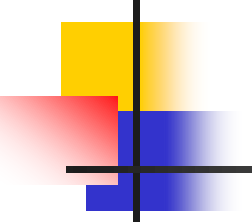
输入向量的每个分量取值离散,
先假设取值 $(0, 1)$, t 对应二类分类

$$\mathbf{x} = [x_1, x_2, \dots, x_D]^T \quad x_i \in \{0, 1\}$$

例 $\mathbf{x} = [0, 1, 1, 0, 1, 0, 0, 1, 0, 0]^T \quad D = 10$

假设 $p(\mathbf{x} | t) = p(x_1, x_2, \dots, x_D | t)$

$$= \prod_{i=1}^D p(x_i | t)$$



离散生成模型: **Naïve Bayes**

用 $t=1$ 表示 C_1 , $t=0$ 表示 C_2

用如下符号

$$\mu_{i|t=1} = \mu_{i|1} = p(x_i = 1|t = 1) = p(x_i = 1|C_1)$$

$$\mu_{i|t=0} = \mu_{i|0} = p(x_i = 1|t = 0) = p(x_i = 1|C_2)$$

类条件概率

$$p(\mathbf{x}|C_k) = \prod_{i=1}^D \mu_{i|k}^{x_i} (1 - \mu_{i|k})^{1-x_i}$$



4.2 Naïve Bayes学习

样本集 $\left\{ \mathbf{x}^{(n)}, t^{(n)} \right\}_{n=1}^N$

注意，区别输入向量下标，用上标表示样本序号

$$p(t = 1) = p(C_1) = \pi$$

$$p(t = 0) = p(C_2) = 1 - \pi$$



Naïve Bayes学习

联合概率表示1

$$\begin{aligned} p(\mathbf{x}^{(n)}, C_1) &= p(\mathbf{x}^{(n)}, t = 1) \\ &= p(C_1)p(\mathbf{x}^{(n)}|C_1) \\ &= \pi \prod_{i=1}^D \mu_{i|1}^{x_i^{(n)}} (1 - \mu_{i|1})^{1-x_i^{(n)}} \end{aligned}$$



Naïve Bayes学习

联合概率表示2

$$\begin{aligned} p(\mathbf{x}^{(n)}, C_2) &= p(\mathbf{x}^{(n)}, t = 0) \\ &= p(C_2)p(\mathbf{x}^{(n)}|C_2) \\ &= (1 - \pi) \prod_{i=1}^D \mu_{i|0}^{x_i^{(n)}} (1 - \mu_{i|0})^{1-x_i^{(n)}} \end{aligned}$$

Naïve Bayes学习

似然函数

$$p(t, X | \pi, \mu_{i|1}, \mu_{i|0})$$

$$= \prod_{n=1}^N \left(\pi \prod_{i=1}^D \mu_{i|1}^{x_i^{(n)}} (1 - \mu_{i|1})^{1-x_i^{(n)}} \right)^{t^{(n)}} \times$$
$$\times \left((1 - \pi) \prod_{i=1}^D \mu_{i|0}^{x_i^{(n)}} (1 - \mu_{i|0})^{1-x_i^{(n)}} \right)^{1-t^{(n)}}$$

Naïve Bayes学习



似然函数

$$\frac{\partial \ln p(t, X | \pi, \mu_{i|1}, \mu_{i|0})}{\partial \pi} = 0$$

得

$$\pi = \frac{1}{N} \sum_{n=1}^N t^{(n)}$$

Naïve Bayes学习

似然函数

$$\frac{\partial \ln p(\mathbf{t}, \mathbf{X} | \pi, \mu_{i|1}, \mu_{i|0})}{\partial \mu_{i|1}} = 0 \quad \text{得}$$

$$\mu_{i|1} = \frac{\sum_{n=1}^N t^{(n)} x_i^{(n)}}{\sum_{n=1}^N t^{(n)}}$$

Naïve Bayes学习

似然函数

$$\frac{\partial \ln p(\mathbf{t}, \mathbf{X} | \pi, \mu_{i|1}, \mu_{i|0})}{\partial \mu_{i|0}} = 0 \quad \text{得}$$

$$\mu_{i|0} = \frac{\sum_{n=1}^N (1 - t^{(n)}) x_i^{(n)}}{\sum_{n=1}^N 1 - t^{(n)}}$$

Naïve Bayes学习

$x_i^{(n)}$ 是二元变量，故参数学习算法重写为

$$\pi = \frac{1}{N} \sum_{n=1}^N \mathbf{1}\{t^{(n)} = 1\}$$

$$\mu_{i|1} = \frac{\sum_{n=1}^N \mathbf{I}\{t^{(n)} = 1 \cap x_i^{(n)} = 1\}}{\sum_{n=1}^N \mathbf{I}\{t^{(n)} = 1\}}$$

$$\mu_{i|0} = \frac{\sum_{n=1}^N \mathbf{I}\{t^{(n)} = 0 \cap x_i^{(n)} = 1\}}{\sum_{n=1}^N \mathbf{I}\{t^{(n)} = 0\}}$$

4.3 Naïve Bayes推断



给出新输入 \mathbf{x} ，进行分类

$$\begin{aligned} p(t = 1|\mathbf{x}) &= p(C_1|\mathbf{x}) \\ &= \frac{p(\mathbf{x}|t = 1)p(t = 1)}{p(\mathbf{x})} \\ &= \frac{\left(\prod_{i=1}^D p(x_i|t = 1)\right)p(t = 1)}{\left(\prod_{i=1}^D p(x_i|t = 1)\right)p(t = 1) + \left(\prod_{i=1}^D p(x_i|t = 0)\right)p(t = 0)} \end{aligned}$$

Naïve Bayes分类

可推广到输入各分量
取值为M个值的情况

带入学习得到的参数，后验类概率为

$$p(t = 1|\mathbf{x})$$

$$= \frac{\pi \prod_{i=1}^D \mu_{i|1}^{x_i} (1 - \mu_{i|1})^{1-x_i}}{\pi \prod_{i=1}^D \mu_{i|1}^{x_i} (1 - \mu_{i|1})^{1-x_i} + (1 - \pi) \prod_{i=1}^D \mu_{i|0}^{x_i} (1 - \mu_{i|0})^{1-x_i}}$$

$$p(t = 0|\mathbf{x}) = 1 - p(t = 1|\mathbf{x})$$

4.4 拉普拉斯平滑克服零概率比值问题

Naïve Bayes分类存在的一个问题

若存在 $\mu_{i|1} = 0$ 和 $\mu_{j|0} = 0$ 则,

$$p(t = 1|\mathbf{x}) = \frac{0}{0}$$

无法做出判断

希望不存在 $\mu_{i|k} = 0$

拉普拉斯平滑 (Laplace smoothing)

设

$$z \in \{1, \dots, k\}$$

且定义

$$\phi_i = p(z = i)$$

有样本集

$$\{z^{(1)}, \dots, z^{(m)}\}$$

标准ML估计:

$$\phi_j = \frac{\sum_{i=1}^m 1\{z^{(i)} = j\}}{m}$$

拉普拉斯平滑为

$$\phi_j = \frac{\sum_{i=1}^m 1\{z^{(i)} = j\} + 1}{m + k}$$

Naïve Bayes参数学习改进为

拉普拉斯平滑参数估计

$$\pi = \frac{1}{N} \sum_{n=1}^N \mathbf{1}\{t^{(n)} = 1\}$$

$$\mu_{i|1} = \frac{\sum_{n=1}^N \mathbf{I}\{t^{(n)} = 1 \cap x_i^{(n)} = 1\} + 1}{\sum_{n=1}^N \mathbf{I}\{t^{(n)} = 1\} + 2}$$

$$\mu_{i|0} = \frac{\sum_{n=1}^N \mathbf{I}\{t^{(n)} = 0 \cap x_i^{(n)} = 1\} + 1}{\sum_{n=1}^N \mathbf{I}\{t^{(n)} = 0\} + 2}$$