

线性分类

频率派 \rightarrow 统计机器学习.

贝叶斯派 \rightarrow 概率图模型.

线性回归
 $f(w, b) = w^T x + b$
 $x \in \mathbb{R}^p$

属性: x 是 p 维的, f 关于 x 是线性的.
全局: $w^T x + b$ 是线性组合, 直接输出.
系数: w 关于 w 是线性的.
全局性: 在整个特征空间拟合一条直线, 全局.
① 数据未加工: 直接使用数据.

线性回归是统计机器学习的基本模型, 其它模型通过打破其中一个或多个特点形成了统计机器学习的架构.

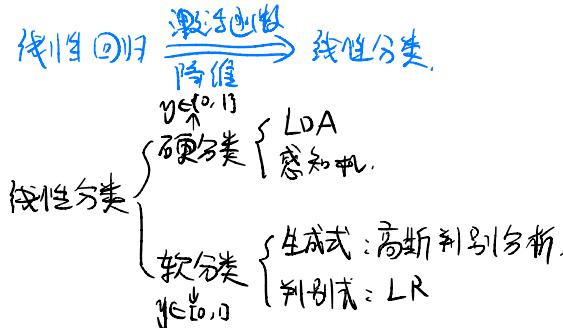
若打破局部线性, 如多项式回归, 其属性可能是多次的, 如¹

全局线性: 如线性分类, 将 x 作为激活函数的输入, 感知机.

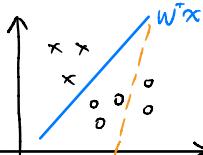
参数非线性: 神经网络, 感知机

若打破全局性: 对样本空间分割, 如纯性决策回归, 决策树.

若打破数据未加工: PCA, 波形.



感知机 (Perceptron)



思想：错误驱动 \rightarrow 指示函数， $\text{sign}(a) = \begin{cases} +1, & a \geq 0 \\ -1, & a < 0 \end{cases}$

模型： $f(x) = \text{sign}(w^T x)$, $x \in \mathbb{R}^p$, $w \in \mathbb{R}^p$, 样本个数为 N .

整体思想：给定初始 w (图中橙线)，会获得被分类错误的样本集合 D ，依次取 D 中的样本，调整 w ，使得 w 对当前取得的样本分类正确。

策略：loss function

$$L(w) = \sum_{i=1}^N I\{\eta_i; w^T x_i < 0\}$$

其中 I 为指示函数，当 I 中条件成立时， $I=1$ ，否则 $I=0$ 。

但当 w 变化 Δw 时， I 可能从 0 变为 1，它是不可导的，因此这个 loss 不合适。

当正确分类时，对任意样本 x_i ，有：

$$\begin{cases} w^T x_i > 0, y_i = +1 \\ w^T x_i < 0, y_i = -1 \end{cases} \Rightarrow y_i w^T x_i > 0.$$

当错误分类时，对任意样本 x_i ，有：

$$\begin{cases} w^T x_i > 0, y_i = -1 \\ w^T x_i < 0, y_i = +1 \end{cases} \Rightarrow y_i w^T x_i < 0.$$

$$L(w) = \sum_{x_i \in D} -y_i w^T x_i$$

算法：SGD (随机梯度下降)

$$w^{(t+1)} \leftarrow w^{(t)} - \lambda \nabla_w L$$

由于每个错误分类的点的 $y_i w^T x_i$ 都小于 0，因此 $-y_i w^T x_i$ 大于 0，此值可作为损失值，将所有误分类点的 $-y_i w^T x_i$ 之和作为损失函数即可。

线性判别分析 (LDA)

$$\text{设: } X = (x_1, x_2, \dots, x_N)^T = \begin{pmatrix} x_1^T \\ x_2^T \\ \vdots \\ x_N^T \end{pmatrix}_{N \times 1}$$

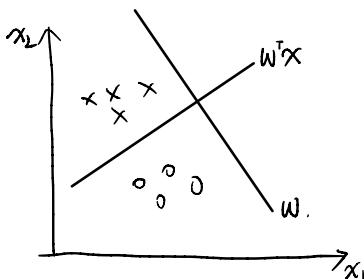
$$Y = \begin{pmatrix} y_1 \\ y_2 \\ \vdots \\ y_N \end{pmatrix}_{N \times 1}$$

$$\{(x_i, y_i)\}_{i=1}^N, x_i \in \mathbb{R}^p, y_i \in \{+1, -1\}$$

$$X_{C_1} = \{x_i \mid y_i = +1\}, X_{C_2} = \{x_i \mid y_i = -1\}$$

$$|X_{C_1}| = N_1, |X_{C_2}| = N_2, N_1 + N_2 = N.$$

LDA 的思想: 类内小, 类间大.



限定 $\|w\|=1$.

样本点在 w 上的投影: $Z_i = w^T x_i$

样本投影的均值: $\bar{Z}_i = \frac{1}{N} \sum_{i=1}^N Z_i = \frac{1}{N} \sum_{i=1}^N w^T x_i$

$$\begin{aligned} \text{Z的协方差矩阵: } S_Z &= \frac{1}{N} \sum_{i=1}^N (Z_i - \bar{Z})(Z_i - \bar{Z})^T \\ &= \frac{1}{N} \sum_{i=1}^N (w^T x_i - \bar{Z})(w^T x_i - \bar{Z})^T \end{aligned}$$

$$C_1: \bar{Z}_1 = \frac{1}{N_1} \sum_{i=1}^{N_1} w^T x_i$$

$$S_1 = \frac{1}{N_1} \sum_{i=1}^{N_1} (w^T x_i - \bar{Z}_1)(w^T x_i - \bar{Z}_1)^T$$

$$C_2: \bar{Z}_2 = \frac{1}{N_2} \sum_{i=1}^{N_2} w^T x_i$$

$$S_2 = \frac{1}{N_2} \sum_{i=1}^{N_2} (w^T x_i - \bar{Z}_2)(w^T x_i - \bar{Z}_2)^T$$

$$\text{类间距离: } (\bar{Z}_1 - \bar{Z}_2)^2$$

$$\text{类内距离: } S_1 + S_2$$

$$\text{目标函数: } J(w) = \frac{(\bar{Z}_1 - \bar{Z}_2)^2}{S_1 + S_2}, \text{ 令 } \hat{w} = \arg \max_w J(w)$$

$$J(w) = \frac{(\bar{Z}_1 - \bar{Z}_2)^2}{S_1 + S_2}$$

$$\text{分子: } \left(\frac{1}{N_1} \sum_{i=1}^{N_1} w^T x_i - \frac{1}{N_2} \sum_{i=1}^{N_2} w^T x_i \right)^2 = \left(w^T \left(\frac{1}{N_1} \sum_{i=1}^{N_1} x_i - \frac{1}{N_2} \sum_{i=1}^{N_2} x_i \right) \right)^2 = \left(w^T (\bar{x}_{C_1} - \bar{x}_{C_2}) \right)^2$$

$$\text{总错} = S_1 + S_2$$

$$\begin{aligned} S_1 &= \frac{1}{N_1} \sum_{i=1}^{N_1} (w^T x_i - \bar{x}_{c_1}) (w^T x_i - \bar{x}_{c_1})^T \\ &= \frac{1}{N_1} \sum_{i=1}^{N_1} w^T (x_i - \bar{x}_{c_1}) (x_i - \bar{x}_{c_1})^T w \\ &= w^T \underbrace{\left[\frac{1}{N_1} \sum_{i=1}^{N_1} (x_i - \bar{x}_{c_1}) (x_i - \bar{x}_{c_1})^T \right]}_{S_{C_1}} w \\ &= w^T S_{C_1} w \end{aligned}$$

$$\text{故 总错} = w^T S_w w + w^T S_{C_2} w.$$

$$= w^T (S_{C_1} + S_{C_2}) w$$

$$\text{故 } J(w) = \frac{w^T (\bar{x}_{c_1} - \bar{x}_{c_2}) (\bar{x}_{c_1} - \bar{x}_{c_2})^T w}{w^T (S_{C_1} + S_{C_2}) w}$$

模型求解: $J(w) = \frac{w^T S_b w}{w^T S_w w}$, 其中 S_b : 类间方差, S_w : 类内方差.

$$J(w) = w^T S_w w (w^T S_w w)^{-1}$$

$$\Leftrightarrow \frac{\partial J(w)}{\partial w} = 2S_b w (w^T S_w w)^{-1} + w^T S_b w + 1 \cdot (w^T S_w w)^{-2} \cdot 2S_w \cdot w = 0.$$

$$\text{即 } S_b w (w^T S_w w) - w^T S_b w S_w w = 0$$

$$\underbrace{w^T S_w w}_{\|x\| \in R} \underbrace{S_b w}_{\|w\| \in R} = S_b w \underbrace{(w^T S_w w)}_{\|w\| \in R}$$

$$\text{即 } S_b w = \frac{w^T S_w w}{w^T S_b w} \cdot S_b w$$

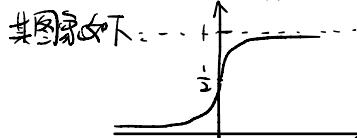
$$\begin{aligned} w &= \frac{w^T S_w w}{w^T S_b w} \cdot S_b^{-1} \cdot S_b \cdot w \propto \underbrace{S_b^{-1} \cdot S_b \cdot w}_{(\bar{x}_{c_1} - \bar{x}_{c_2}) (\bar{x}_{c_1} - \bar{x}_{c_2})^T w} \\ &\propto S_w^{-1} \cdot (\bar{x}_{c_1} - \bar{x}_{c_2}) \end{aligned}$$

当类另1有K个时, 最后降维到 $K-1$ 维.

Logistic Regression.

判别模型直接对 $P(Y|X)$ 进行建模，应用 MLE.

$$\text{sigmoid 函数: } \sigma(z) = \frac{1}{1+e^{-z}}$$



$z \rightarrow +\infty, \sigma(z) \rightarrow 1$

$z = 0, \sigma(z) = \frac{1}{2}$

$z \rightarrow -\infty, \sigma(z) \rightarrow 0$

$$\begin{cases} p_1 = P(y=1|x) = \sigma(w^T x) = \frac{1}{1+e^{-w^T x}}, y=1 \\ p_0 = P(y=0|x) = 1 - \frac{1}{1+e^{-w^T x}} = \frac{e^{w^T x}}{1+e^{w^T x}}, y=0. \end{cases}$$

$\Rightarrow P(y|x) = p_1^y p_0^{1-y}$

$$\begin{aligned} \text{MLE: } \hat{w} &= \arg \max_w \log P(Y|X) \\ &= \arg \max_w \log \prod_{i=1}^n P(y_i|x_i) \\ &= \arg \max_w \sum_{i=1}^n \log P(y_i|x_i) \\ &= \arg \max_w \underbrace{\sum_{i=1}^n (y_i \log p_1 + (1-y_i) \log p_0)}_{-\text{cross entropy}} \end{aligned}$$

MLE $\xrightarrow{\text{(max)}}$ loss function (min cross entropy)

高斯判别分析 (Gaussian Discriminant Analysis)

Data: $\{(x_i, y_i)\}_{i=1}^N$, $x_i \in \mathbb{R}^p$, $y_i \in \{0, 1\}$

前面的线性分类模型都是判别模型，直接对 $P(Y|X)$ 进行建模，而 $\hat{y} = \arg \max_{y \in \{0, 1\}} P(y|X)$
而高斯判别模型是生成模型，应用贝叶斯定理: $P(Y|X) = \frac{P(X|Y) \cdot P(Y)}{P(X)} \propto P(X|Y) \cdot P(Y)$
 $= P(X|Y)$, 所以生成模型是对联合概率进行建模。

即求 $\hat{y} = \arg \max_{y \in \{0, 1\}} P(y) \cdot P(X|y)$

GDA:

对于 y , 可将 y 看作伯努利分布, $y \sim \text{Bernoulli}(\phi)$

$$\left. \begin{aligned} x|y=1 &\sim N(\mu_1, \Sigma) \\ x|y=0 &\sim N(\mu_2, \Sigma) \end{aligned} \right\} \Rightarrow N(\mu_1, \Sigma)^y \cdot N(\mu_2, \Sigma)^{1-y}$$

y	1	0
	ϕ	$1-\phi$

$$P(y) = \phi^y \cdot (1-\phi)^{1-y}$$

$$\begin{aligned} \text{log-likelihood} &= l(\theta) = \log \prod_{i=1}^N P(x_i, y_i) \\ \theta = (\mu_1, \mu_2, \Sigma, \phi) &= \sum_{i=1}^N \log [P(x_i|y_i) \cdot P(y_i)] \\ &= \sum_{i=1}^N [\log P(x_i|y_i) + \log P(y_i)] \\ &= \sum_{i=1}^N [\log N(\mu_1, \Sigma)^{y_i} \cdot N(\mu_2, \Sigma)^{1-y_i} + \log \phi^{y_i} \cdot (1-\phi)^{1-y_i}] \\ &= \sum_{i=1}^N \underbrace{\log N(\mu_1, \Sigma)^{y_i}}_{①} + \underbrace{\log N(\mu_2, \Sigma)^{1-y_i}}_{②} + \underbrace{\log \phi^{y_i} \cdot (1-\phi)^{1-y_i}}_{③} \\ \hat{\theta} &= \arg \max_{\theta} l(\theta) \end{aligned}$$

模型求解

求 ϕ :

$$③ = \sum_{i=1}^N y_i \log \phi + (1-y_i) \log (1-\phi)$$

$$\frac{\partial ③}{\partial \phi} = \sum_{i=1}^N y_i \cdot \frac{1}{\phi} - (1-y_i) \cdot \frac{1}{1-\phi} = 0$$

$$\Rightarrow \sum_{i=1}^N y_i \cdot (1-\phi) - (1-y_i) \phi = 0$$

$$\sum_{i=1}^N y_i - y_i \phi - \phi + y_i \phi = 0$$

$$\sum_{i=1}^N (y_i - \phi) = 0$$

$$\sum_{i=1}^N y_i - N\phi = 0$$

$$\therefore \boxed{\hat{\phi} = \frac{1}{N} \sum_{i=1}^N y_i}$$

设 $y=1$ 的样本为 N_1 个, $y=0$ 的样本 N_2 个, $N=N_1+N_2$

$$\text{则 } \hat{\mu} = \frac{1}{N} \sum_{i=1}^N y_i = \frac{N_1}{N}$$

求 μ_1, μ_2 :

$$\textcircled{1} = \sum_{i=1}^N \log N(\mu_1, \Sigma)^{y_i}$$

$$= \sum_{i=1}^N y_i \log \frac{1}{(2\pi)^{\frac{D}{2}} |\Sigma|^{\frac{1}{2}}} \exp \left(-\frac{1}{2} (x_i - \mu)^T \Sigma^{-1} (x_i - \mu) \right)$$

$$\mu_1 = \arg \max_{\mu_1} \textcircled{1} = \arg \max_{\mu_1} \sum_{i=1}^N y_i \left(-\frac{1}{2} (x_i - \mu)^T \Sigma^{-1} (x_i - \mu) \right)$$

$$\Delta = \sum_{i=1}^N y_i \left(-\frac{1}{2} (x_i - \mu)^T \Sigma^{-1} (x_i - \mu) \right) = -\frac{1}{2} \sum_{i=1}^N y_i (x_i^T \Sigma^{-1} x_i - x_i^T \Sigma^{-1} \mu_1 - \mu_1^T \Sigma^{-1} x_i + \mu_1^T \Sigma^{-1} \mu_1)$$

$$= -\frac{1}{2} \sum_{i=1}^N y_i (x_i^T \Sigma^{-1} x_i - 2\mu_1^T \Sigma^{-1} x_i + \mu_1^T \Sigma^{-1} \mu_1)$$

$$\frac{\partial \Delta}{\partial \mu_1} = -\frac{1}{2} \sum_{i=1}^N y_i (-2 \sum_i x_i + 2 \sum_i \mu_1) = 0$$

$$\sum_{i=1}^N y_i (\Sigma^{-1} \mu_1 - \sum_i x_i) = 0$$

$$\sum_{i=1}^N y_i (\mu_1 - x_i) = 0$$

$$\boxed{R(\hat{\mu}_1) = \frac{\sum_{i=1}^N y_i x_i}{\sum_{i=1}^N y_i} = \frac{\sum_{i=1}^N y_i x_i}{N_1}}$$

μ_2 同理.

求 Σ :

设 $y_i=1$ 的样本点集合为 C_1 , $y_i=0$ 的样本点集合为 C_2

$$\textcircled{2} = \arg \max_{\Sigma} \textcircled{1} + \textcircled{2}$$

$$\textcircled{1} + \textcircled{2} = \sum_{x_i \in C_1} \log N(\mu_1, \Sigma) + \sum_{x_i \in C_2} \log N(\mu_2, \Sigma)$$

$$\begin{aligned} \text{由 } \sum_{i=1}^N \log N(\mu, \Sigma) &= \sum_{i=1}^N \left[\log \frac{1}{(2\pi)^{\frac{D}{2}} |\Sigma|^{\frac{1}{2}}} \exp \left\{ -\frac{1}{2} (x_i - \mu)^T \Sigma^{-1} (x_i - \mu) \right\} \right] \\ &= \sum_{i=1}^N \left[\log \frac{1}{(2\pi)^{\frac{D}{2}}} + \log |\Sigma|^{-\frac{1}{2}} + \left(-\frac{1}{2} (x_i - \mu)^T \Sigma^{-1} (x_i - \mu) \right) \right] \\ &= \sum_{i=1}^N \left[C - \frac{1}{2} \log |\Sigma| - \frac{1}{2} (x_i - \mu)^T \Sigma^{-1} (x_i - \mu) \right] \\ &= C - \frac{1}{2} N \log |\Sigma| - \underbrace{\frac{1}{2} \sum_{i=1}^N (x_i - \mu)^T \Sigma^{-1} (x_i - \mu)}_{\text{实数的和等于它本身.}} \end{aligned}$$

$$\begin{aligned} &\stackrel{\text{对称性}}{=} \sum_{i=1}^N \text{tr}((x_i - \mu)^T \Sigma^{-1} (x_i - \mu)) \\ &= \sum_{i=1}^N \text{tr}((x_i - \mu)(x_i - \mu)^T \Sigma^{-1}) \quad \text{对称矩阵.} \\ &= \text{tr} \left(\sum_{i=1}^N (x_i - \mu)(x_i - \mu)^T \Sigma^{-1} \right) \end{aligned}$$

$$= -\frac{1}{2}N \cdot \log |\Sigma| - \frac{1}{2}N \cdot \text{tr}(S \cdot \Sigma^{-1}) + C$$

$$\text{④ } ①+② = -\frac{1}{2}N_1 \log |\Sigma| - \frac{1}{2}N_1 \cdot \text{tr}(S_1 \cdot \Sigma^{-1}) - \frac{1}{2}N_2 \log |\Sigma| - \frac{1}{2}N_2 \cdot \text{tr}(S_2 \cdot \Sigma^{-1}) + C$$

$$= -\frac{1}{2}[N \log |\Sigma| + N_1 \cdot \text{tr}(S_1 \cdot \Sigma^{-1}) + N_2 \cdot \text{tr}(S_2 \cdot \Sigma^{-1})] + C$$

$$\frac{\partial ④+②}{\partial \Sigma} = -\frac{1}{2}(N\Sigma^{-1} + N_1 \cdot \frac{\partial \text{tr}(S_1 \cdot \Sigma^{-1})}{\partial \Sigma^{-1}} \cdot \frac{\partial \Sigma^{-1}}{\partial \Sigma} + N_2 \cdot \frac{\partial \text{tr}(S_2 \cdot \Sigma^{-1})}{\partial \Sigma^{-1}} \cdot \frac{\partial \Sigma^{-1}}{\partial \Sigma})$$

$$= -\frac{1}{2}(N \cdot \Sigma^{-1} - N_1 S_1^\top \Sigma^{-2} - N_2 S_2^\top \Sigma^{-2})$$

$$= -\frac{1}{2}(N \cdot \Sigma^{-1} - N_1 S_1^\top \Sigma^{-2} - N_2 S_2^\top \Sigma^{-2}) = 0$$

$$\text{⑤ } N\Sigma - N_1 S_1 - N_2 S_2 = 0$$

$$\boxed{\hat{\Sigma} = \frac{N_1 S_1 + N_2 S_2}{N}}$$