

Section 4

Linear Classification

- ① 背景
- ② 感知机
- ③ 线性判别分析
- ④ 逻辑回归
- ⑤ 高斯判别分析
- ⑥ 朴素贝叶斯分类器

背景 (Background)

{ 频率派 → 统计机器学习 核心是线性回归.
贝叶斯派 → 概率图模型

Linear Regression

① 线性

属性线性
系数线性
结果线性

$f(w, b) = w^T x + b$ ^{可以合并到 w 中}
 f 关于 x 是线性的, 关于 w 和 b 也是线性的
→ 直接将线性组合作为结果

② 全局性

对数据的拟合是全局的, 不会局部拟合

③ 数据未加工

数据没有经过降维等处理.

将线性回归这三个性质分别破坏后, 就形成了统计机器学习的大框架.

① 线性 \Rightarrow 非线性

属性非线性

特征转换 \rightarrow 多项式回归

$w_1 x_1^2 + w_2 x_2^2 + \dots$

系数非线性

神经网络. 这里的系数非线性是指系数会变化

系数其实不存在非线性

w^2 也只是系数, 可以令 $\beta = w^2$ 保持线性

神经网络初始化参数不同, 最终得到系数也不同.

结果非线性

线性分类 $\rightarrow f(w^T x + b)$

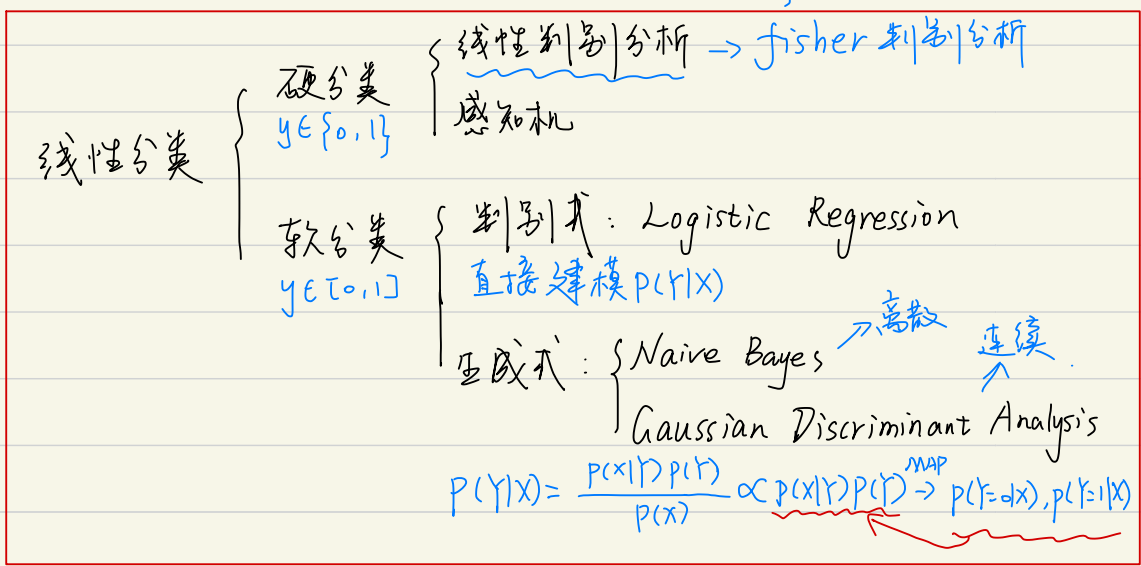
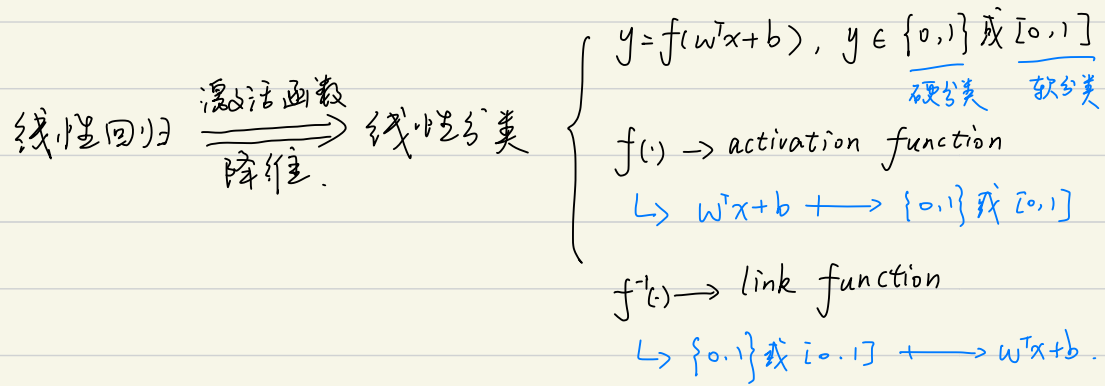
$f(\cdot)$ 是非线性函数, 使得结果输出为非线性.

② 全局性 \Rightarrow 局部性

线性样条回归. 将数据分段拟合, 具有局部性.

决策树. 将样本空间划分为一个个子空间, 使其局部具有不同决策.

③ 数据未加工 \Rightarrow 处理数据. PCA. 流形



感知机 (Perceptron)

思想: 错误驱动

模型: $f(x) = \text{sign}(w^T x)$, $x \in \mathbb{R}^p$, $w \in \mathbb{R}^p$, $\text{sign}(a) = \begin{cases} +1, & a \geq 0 \\ -1, & a < 0 \end{cases}$

策略: Loss Function 被错误分类的样本的数量.

$$L'(w) = \sum_{i=1}^N \mathbb{I}\{y_i w^T x_i < 0\}$$

错误分类时, y_i 与 $w^T x_i$ 异号.

但是 $w \rightarrow w + \Delta w$, \mathbb{I} 就会从 0 阶跃到 1 或从 1 阶跃到 0

即 $L'(w)$ 不可导.

\therefore 不用 $L'(w)$ 作为损失函数.

令 $D = \{\text{误分类的样本}\}$

则 $L(w) = \sum_{x_i \in D} -y_i w^T x_i$, 将其作为损失函数.

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

参数的更新采用梯度下降法: 全体样本

$$w^{(t+1)} \leftarrow w^{(t)} - \lambda \left(\sum_{x_i \in D} -y_i x_i \right) = w^{(t)} + \lambda \sum_{x_i \in D} y_i x_i$$

或随机梯度下降: 单个样本.

$$w^{(t+1)} \leftarrow w^{(t)} - \lambda (-y_i x_i) = w^{(t)} + \lambda y_i x_i$$

线性可分的数据用感知机算法才能收敛.

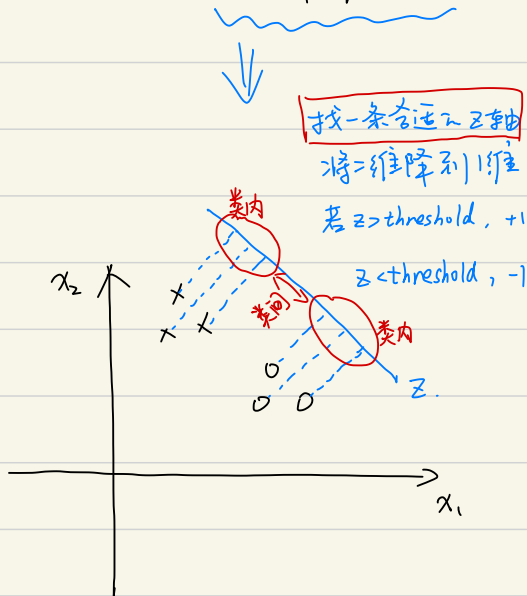
若线性不可分, 则需要使用感知机的变形 pocket algorithm.

允许有错误率

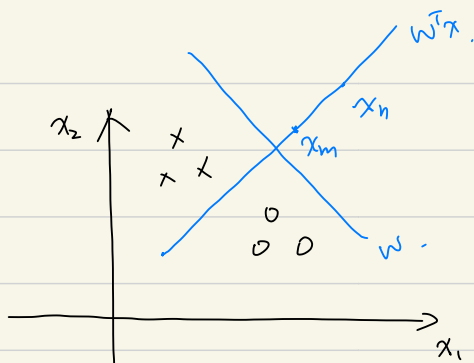
线性判别分析 (Linear Discriminant Analysis)

设 $X = (x_1, x_2, \dots, x_N)^T \in \mathbb{R}^{N \times P}$
 $Y = (y_1, y_2, \dots, y_N)^T \in \mathbb{R}^{N \times 1}$, $y_i \in \left\{ \overset{c_1}{\uparrow} +1, \overset{c_2}{\uparrow} -1 \right\}$
 $X_{c_1} = \{x_i | y_i = +1\}$, X_{c_1} 集合大小为 N_1
 $X_{c_2} = \{x_i | y_i = -1\}$, X_{c_2} 集合大小为 N_2
 $N_1 + N_2 = N$

思想：类内小，类间大，运用降维的方式。



假设超平面 $w^T x = 0$ 将两类样本分割开，
 可以证明，投影的轴即为超平面的法向量 w 。



证明：由于 $w^T x = 0$ 。

任取超平面上两点

$$\text{得 } w^T x_m = 0, w^T x_n = 0$$

$$\text{则 } w^T (x_m - x_n) = 0.$$

由于 w 与 $(x_m - x_n)$ 内积为 0。则 w 垂直于 $(x_m - x_n)$

由于 x_m, x_n 是任意选取的，则 w 为超平面 $w^T x = 0$ 的法向量

规定超平面 $w^T x = 0$ ， $\|w\| = 1$ 。

则样本点 x_i 在超平面法向量上的投影长度为 $z_i = w^T x_i$ 内积

$$\text{则均值 } \bar{z} = \frac{1}{N} \sum_{i=1}^N z_i = \frac{1}{N} \sum_{i=1}^N w^T x_i$$

$$\text{方差 } S = \frac{1}{N} \sum_{i=1}^N (z_i - \bar{z})(z_i - \bar{z})^T$$

$$\text{对于 } C_1 \text{ 类投影，均值 } \bar{z}_1 = \frac{1}{N_1} \sum_{x_i \in X_{C_1}} w^T x_i$$

$$\text{方差 } S_{z_1} = \frac{1}{N_1} \sum_{x_i \in X_{C_1}} (w^T x_i - \bar{z}_1)(w^T x_i - \bar{z}_1)^T$$

$$\text{对于 } C_2 \text{ 类投影，均值 } \bar{z}_2 = \frac{1}{N_2} \sum_{x_i \in X_{C_2}} w^T x_i$$

$$\text{方差 } S_{z_2} = \frac{1}{N_2} \sum_{x_i \in X_{C_2}} (w^T x_i - \bar{z}_2)(w^T x_i - \bar{z}_2)^T$$

这里需要注意的是， z_i 为投影后的向量长度。

同时, 对于C1类本身 $\bar{x}_{c1} = \frac{1}{N_1} \sum_{x_i \in X_{c1}} x_i$

$$\text{方差: } S_{c1} = \frac{1}{N_1} \sum_{x_i \in X_{c1}} (x_i - \bar{x}_{c1})(x_i - \bar{x}_{c1})^T$$

对于C2类本身: $\bar{x}_{c2} = \frac{1}{N_2} \sum_{x_i \in X_{c2}} x_i$

$$\text{方差: } S_{c2} = \frac{1}{N_2} \sum_{x_i \in X_{c2}} (x_i - \bar{x}_{c2})(x_i - \bar{x}_{c2})^T$$

则类间方差: $S_b = (\bar{x}_{c1} - \bar{x}_{c2})(\bar{x}_{c1} - \bar{x}_{c2})^T$ between-class

类内方差: $S_w = S_{c1} + S_{c2}$ within-class

注意: 这里的 \bar{x}_{c1} , S_{c1} , \bar{x}_{c2} , S_{c2} , S_b , S_w 都不是标量。
是针对原始空间来说的,

为了使投影后实现“类内小, 类间大”, 设目标函数:

$$J(w) = \frac{(\bar{z}_1 - \bar{z}_2)^2}{S_{z1} + S_{z2}} \quad \text{则 } \hat{w} = \arg \max_w J(w)$$

$$\begin{aligned} \text{对于分子: } (\bar{z}_1 - \bar{z}_2)^2 &= (\bar{z}_1 - \bar{z}_2)(\bar{z}_1 - \bar{z}_2)^T \\ &= \left(\frac{1}{N_1} \sum_{x_i \in X_{c1}} w^T x_i - \frac{1}{N_2} \sum_{x_i \in X_{c2}} w^T x_i \right) \left(\frac{1}{N_1} \sum_{x_i \in X_{c1}} w^T x_i - \frac{1}{N_2} \sum_{x_i \in X_{c2}} w^T x_i \right)^T \\ &= w^T \left(\frac{1}{N_1} \sum_{x_i \in X_{c1}} x_i - \frac{1}{N_2} \sum_{x_i \in X_{c2}} x_i \right) \left(\frac{1}{N_1} \sum_{x_i \in X_{c1}} x_i - \frac{1}{N_2} \sum_{x_i \in X_{c2}} x_i \right)^T w \\ &= w^T (\bar{x}_{c1} - \bar{x}_{c2})(\bar{x}_{c1} - \bar{x}_{c2})^T w \\ &= w^T S_b w. \end{aligned}$$

$$\begin{aligned}
\text{对于分母: } S_{z_1} + S_{z_2} &= \frac{1}{N_1} \sum_{x_i \in X_{c_1}} (w^T x_i - \frac{1}{N_1} \sum_{x_i \in X_{c_1}} w^T x_i) (w^T x_i - \frac{1}{N_1} \sum_{x_i \in X_{c_1}} w^T x_i)^T \\
&\quad + \frac{1}{N_2} \sum_{x_i \in X_{c_2}} (w^T x_i - \frac{1}{N_2} \sum_{x_i \in X_{c_2}} w^T x_i) (w^T x_i - \frac{1}{N_2} \sum_{x_i \in X_{c_2}} w^T x_i)^T \\
&= \frac{1}{N_1} \sum_{x_i \in X_{c_1}} w^T (x_i - \bar{x}_{c_1}) (x_i - \bar{x}_{c_1})^T w \\
&\quad + \frac{1}{N_2} \sum_{x_i \in X_{c_2}} w^T (x_i - \bar{x}_{c_2}) (x_i - \bar{x}_{c_2})^T w \\
&= w^T \left[\frac{1}{N_1} \sum_{x_i \in X_{c_1}} (x_i - \bar{x}_{c_1}) (x_i - \bar{x}_{c_1})^T \right] w \\
&\quad + w^T \left[\frac{1}{N_2} \sum_{x_i \in X_{c_2}} (x_i - \bar{x}_{c_2}) (x_i - \bar{x}_{c_2})^T \right] w \\
&= w^T S_{c_1} w + w^T S_{c_2} w \\
&= w^T (S_{c_1} + S_{c_2}) w \\
&= w^T S_w w
\end{aligned}$$

$$\text{所以目标函数: } J(w) = \frac{w^T S_b w}{w^T S_w w}$$

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

$$= 0$$

$$\lambda | 2S_b W (W^T S_w W)^{-1} = W^T S_b W (W^T S_w W)^{-2} \geq S_w W \quad \text{都为实数}$$

$$S_b W (W^T S_w W) = W^T S_b W S_w W \leftarrow \text{两边同乘}(W^T S_w W)^2$$

$$\frac{W^T S_w W}{W^T S_b W} S_b W = S_w W$$

$$W = \boxed{\frac{W^T S_w W}{W^T S_b W}} S_w^{-1} S_b W \quad \text{一维}$$

由于 $W^T S_w W$ 和 $W^T S_b W$ 都是一维实数，则

由于 W 只考虑其方向，大小不用考虑（前面约束 $\|W\|=1$ ，就算前面没有约束，最后可以 scaling 成 1），
所以可以直接舍去一维实数。

$$W \propto S_w^{-1} S_b W$$

$$\propto S_w^{-1} (\bar{X}_{C_1} - \bar{X}_{C_2}) (\bar{X}_{C_1} - \bar{X}_{C_2})^T W \quad \text{一维}$$

$$\propto S_w^{-1} (\bar{X}_{C_1} - \bar{X}_{C_2})$$

至此，求得 W 的方向。进一步，如果 C_1 类与 C_2 类 各向同性。

$$\text{那么 } S_w = \lambda I. \text{ 即 } W \propto (\bar{X}_{C_1} - \bar{X}_{C_2})$$

协方差矩阵为对角

矩阵，且对角线元素相同

综上，求得投影的方向，满足“类内小，类间大”的要求。

逻辑回归 (Logistic Regression)

线性回归的结果是线性组合: $w^T x$

逻辑回归的结果是类别: $\{0, 1\}$.

用激活函数将 $w^T x$ 映射到 $\{0, 1\}$.

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

$$\lim_{z \rightarrow \infty} \sigma(z) = 1$$

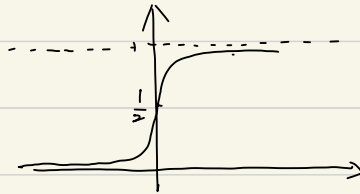
$$z \rightarrow \infty$$

$$\lim_{z \rightarrow -\infty} \sigma(z) = 0$$

$$z \rightarrow -\infty$$

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

$$z \rightarrow 0$$



$$\text{知} \mid \sigma: \mathbb{R} \rightarrow (0, 1)$$

$$w^T x \mapsto p$$

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

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

用极大似然估计就可以直接估计 w .

$$MLE: \hat{w} = \underset{w}{\operatorname{argmax}} \log P(Y|X)$$

$$= \underset{w}{\operatorname{argmax}} \log \prod_{i=1}^N P(y_i | x_i)$$

$$= \underset{w}{\operatorname{argmax}} \sum_{i=1}^N \log P(y_i | x_i)$$

$$= \underset{w}{\operatorname{argmax}} \sum_{i=1}^N [y_i \log p_1 + (1-y_i) \log p_0]$$

$$= \underset{w}{\operatorname{argmax}} \sum_{i=1}^N [y_i \log P(y=1|x_i) + (1-y_i) \log (1 - P(y=1|x_i))]$$

$(-1) \times$ cross entropy

$$= \underset{w}{\operatorname{argmax}} \sum_{i=1}^N [-y_i \log (1 + e^{-w^T x_i}) - (1-y_i) (w^T x_i + \log (1 + e^{-w^T x_i}))]$$

max MLE \Rightarrow min Loss function

$$= \underset{w}{\operatorname{argmin}} \sum_{i=1}^N [y_i (\log (1 + e^{w^T x_i}) - w^T x_i) + (1-y_i) \log (1 + e^{w^T x_i})]$$

$L(w)$

$$\frac{\partial L(w)}{\partial w} = \sum_{i=1}^N \left[y_i \left(\frac{x_i e^{w^T x_i}}{1 + e^{w^T x_i}} - x_i \right) + (1-y_i) \left(-\frac{x_i e^{w^T x_i}}{1 + e^{w^T x_i}} \right) \right]$$

$$= \sum_{i=1}^N \left[y_i \frac{x_i e^{w^T x_i}}{1 + e^{w^T x_i}} - y_i x_i + \frac{x_i e^{w^T x_i}}{1 + e^{w^T x_i}} - y_i \frac{x_i e^{w^T x_i}}{1 + e^{w^T x_i}} \right]$$

$$= \sum_{i=1}^N \left(\frac{x_i e^{w^T x_i}}{1 + e^{w^T x_i}} - y_i x_i \right)$$

2. 使用梯度下降法, 即可求得 w 近似解:

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

Δ
步长

使用牛顿法也可以. 只不过需要求二阶导数.

高斯判别分析

(Gaussian Discriminant Analysis)

Data: $\{(x_i, y_i)\}$ $x_i \in \mathbb{R}^P$, $y_i \in \{0, 1\}$

高斯判别分析属于前面提到的生成式模型。

对于判别式模型，是求 $P(y|x)$ 的大小。

对 $P(y|x)$ 建模

对于生成式模型，是求 $P(y=0|x)$ 和 $P(y=1|x)$ 哪个大

对 $P(x, y)$ 建模

因此，根据贝叶斯公式： $P(y|x) = \frac{P(x|y)P(y)}{P(x)}$

可得 $P(y|x) \propto P(x|y)P(y)$ 只需要考虑这个联合分布即可， $P(x, y)$

$$\text{MAP: } \hat{y} = \arg \max_{y \in \{0, 1\}} P(y|x) = \arg \max_y P(y) \cdot P(x|y)$$

由于 y 只有两个取值，故此时可以看作是伯努利分布

$$y \sim \text{Bernoulli}(\phi) \quad \begin{array}{c|c} y & 0 & 1 \\ \hline P & 1-\phi & \phi \end{array} \Rightarrow P(y) = \phi^y (1-\phi)^{1-y}$$

$$x|y=1 \sim N(\mu_1, \Sigma)$$

$$x|y=0 \sim N(\mu_2, \Sigma)$$

高斯判别分析的假设

$$\Rightarrow x|y \sim \underbrace{N(\mu_1, \Sigma)}_{P_1}^y \cdot \underbrace{N(\mu_2, \Sigma)}_{P_2}^{1-y} = P_1^y P_2^{1-y}$$

联合概率的似然函数: $L(\theta) = \log \prod_{i=1}^N p(x_i, y_i)$

$$\begin{aligned}
 &= \sum_{i=1}^N \log p(x_i | y_i) p(y_i) \\
 &= \sum_{i=1}^N [\log p(x_i | y_i) + \log p(y_i)] \\
 &= \sum_{i=1}^N [\log p_1^{y_i} p_2^{1-y_i} + \log \phi^{y_i} (1-\phi)^{1-y_i}] \\
 &= \sum_{i=1}^N [y_i \log p_1 + (1-y_i) \log p_2 + \log \phi^{y_i} (1-\phi)^{1-y_i}]
 \end{aligned}$$

其中, $\theta = (\phi, \mu_1, \mu_2, \Sigma)$

2.1 MLE: $\theta = \arg \max_{\theta} L(\theta)$

对于 ϕ :

$$\begin{aligned}
 \frac{\partial L(\theta)}{\partial \phi} &= \frac{\partial}{\partial \phi} \left[\sum_{i=1}^N \log \phi^{y_i} (1-\phi)^{1-y_i} \right] \\
 &= \sum_{i=1}^N [y_i \log \phi + (1-y_i) \log (1-\phi)] \\
 &= \sum_{i=1}^N \left[\frac{y_i}{\phi} - \frac{1-y_i}{1-\phi} \right]
 \end{aligned}$$

令 $\frac{\partial L(\theta)}{\partial \phi} = 0$.

2.1 $\sum_{i=1}^N [y_i - y_i \phi - (1-y_i) \phi] = 0$

$$\phi = \frac{1}{N} \left(\sum_{i=1}^N y_i \right) = \frac{N_1}{N}$$

$$\text{对于 } \mu_1, \frac{\partial L(\theta)}{\partial \mu_1} = \frac{\partial}{\partial \mu_1} \left[\sum_{i=1}^N y_i \log p_i \right]$$

$$\text{其中 } p_i = \frac{1}{(2\pi)^{\frac{1}{2}} |\Sigma|^{\frac{1}{2}}} \exp \left\{ -\frac{1}{2} (x - \mu_1)^T \Sigma^{-1} (x - \mu_1) \right\}$$

$$\begin{aligned} \text{所以 } \frac{\partial L(\theta)}{\partial \mu_1} &= \frac{\partial}{\partial \mu_1} \sum_{i=1}^N -\frac{1}{2} y_i (x_i - \mu_1)^T \Sigma^{-1} (x_i - \mu_1) \\ &= \frac{\partial}{\partial \mu_1} \sum_{i=1}^N -\frac{1}{2} y_i (x_i^T \Sigma^{-1} - \mu_1^T \Sigma^{-1}) (x_i - \mu_1) \\ &= \frac{\partial}{\partial \mu_1} \sum_{i=1}^N -\frac{1}{2} y_i [x_i^T \Sigma^{-1} x_i - 2 x_i^T \Sigma^{-1} \mu_1 + \mu_1^T \Sigma^{-1} \mu_1] \\ &= \sum_{i=1}^N -\frac{1}{2} y_i (-2 \Sigma^{-1} x_i + 2 \Sigma^{-1} \mu_1) \end{aligned}$$

$$\frac{\partial L(\theta)}{\partial \mu_1} = 0$$

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

$$\mu_1 \sum_{i=1}^N y_i = \sum_{i=1}^N y_i x_i$$

$$\mu_1 = \frac{\sum_{i=1}^N y_i x_i}{\sum_{i=1}^N y_i} \quad N_1$$

同理: 对于 μ_2 : $\mu_2 = \frac{\sum_{i=1}^N (1-y_i) x_i}{\sum_{i=1}^N (1-y_i)} \quad N_2$

$N_1 + N_2 = N$

$$\begin{aligned}
\text{对于 } \Sigma: \frac{\partial \mathcal{L}(\theta)}{\partial \Sigma} &= \frac{\partial}{\partial \Sigma} \left[\sum_{i=1}^N y_i \log \frac{1}{(2\pi)^{\frac{p}{2}} |\Sigma|^{\frac{1}{2}}} \exp \left\{ -\frac{1}{2} (x_i - \mu_1)^T \Sigma^{-1} (x_i - \mu_1) \right\} \right. \\
&\quad \left. + \sum_{i=1}^N (1-y_i) \log \frac{1}{(2\pi)^{\frac{p}{2}} |\Sigma|^{\frac{1}{2}}} \exp \left\{ -\frac{1}{2} (x_i - \mu_2)^T \Sigma^{-1} (x_i - \mu_2) \right\} \right] \\
&= \frac{\partial}{\partial \Sigma} \left[\sum_{i=1}^N y_i \left(-\frac{1}{2} \log |\Sigma| - \frac{1}{2} (x_i - \mu_1)^T \Sigma^{-1} (x_i - \mu_1) \right) \right. \\
&\quad \left. + \sum_{i=1}^N (1-y_i) \left(-\frac{1}{2} \log |\Sigma| - \frac{1}{2} (x_i - \mu_2)^T \Sigma^{-1} (x_i - \mu_2) \right) \right] \\
&= \frac{\partial}{\partial \Sigma} \left[\sum_{i=1}^N -\frac{1}{2} \log |\Sigma| + \sum_{i=1}^N -\frac{1}{2} y_i (x_i - \mu_1)^T \Sigma^{-1} (x_i - \mu_1) \right. \\
&\quad \left. + \sum_{i=1}^N -\frac{1}{2} (1-y_i) (x_i - \mu_2)^T \Sigma^{-1} (x_i - \mu_2) \right].
\end{aligned}$$

$$\triangleq C_1 = \{x_i \mid y_i = 1, i = 1, 2, \dots, N\}$$

$$C_2 = \{x_i \mid y_i = 0, i = 1, 2, \dots, N\}.$$

且 C_1 中元素个数为 N_1 , C_2 中为 N_2 . $N_1 + N_2 = N$.

$$\begin{aligned}
\text{则} \quad \frac{\partial \mathcal{L}(\theta)}{\partial \Sigma} &= \frac{\partial}{\partial \Sigma} \left[-\frac{1}{2} N \log |\Sigma| - \frac{1}{2} \sum_{x_i \in C_1} (x_i - \mu_1)^T \Sigma^{-1} (x_i - \mu_1) \right. \\
&\quad \left. - \frac{1}{2} \sum_{x_i \in C_2} (x_i - \mu_2)^T \Sigma^{-1} (x_i - \mu_2) \right]
\end{aligned}$$

$$= \underbrace{-\frac{1}{2} N \cdot \Sigma^{-1}} + \frac{1}{2} \sum_{x_i \in C_1} (x_i - \mu_1) (x_i - \mu_1)^T \cdot \Sigma^{-2}$$

$$+ \frac{1}{2} \sum_{x_i \in C_2} (x_i - \mu_2) (x_i - \mu_2)^T \Sigma^{-2}$$

$$\frac{\partial |A|}{\partial A} = |A| A^{-1}$$

$$\frac{\partial X^T A X}{\partial A} = X X^T$$

$$\frac{\partial A^{-1}}{\partial A} = -A^{-2}$$

其中 $S_1 = \frac{1}{N_1} \sum_{x_i \in C_1} (x_i - \mu_1)(x_i - \mu_1)^T$ C_1 的样本方差

$S_2 = \frac{1}{N_2} \sum_{x_i \in C_2} (x_i - \mu_2)(x_i - \mu_2)^T$ C_2 的样本方差

k. $\left| \frac{\partial L(\theta)}{\partial \Sigma} = -\frac{1}{2} N \cdot \Sigma^{-1} + \frac{1}{2} N_1 S_1 \cdot \Sigma^{-2} + \frac{1}{2} N_2 S_2 \cdot \Sigma^{-2} = 0. \right.$

k. $\left| N \Sigma = N_1 S_1 + N_2 S_2 \right.$

$\Sigma = \frac{1}{N} (N_1 S_1 + N_2 S_2).$

原白板系列中用到的方法和我写的这个稍微有点不同.

原视频中, 将 $(x - \mu)^T \Sigma^{-1} (x - \mu)$ 看作是 $\text{tr}((x - \mu)^T \Sigma^{-1} (x - \mu))$

一个数的迹是它本身.

并利用 $\left| \text{tr}(AB) = \text{tr}(BA) \right|$ 三个公式求它的导数

$\left| \text{tr}(ABC) = \text{tr}(CAB) = \text{tr}(BCA) \right|$

$\left| \frac{\partial \text{tr}(AB)}{\partial A} = B^T \right|$

其实实质上还是一样的

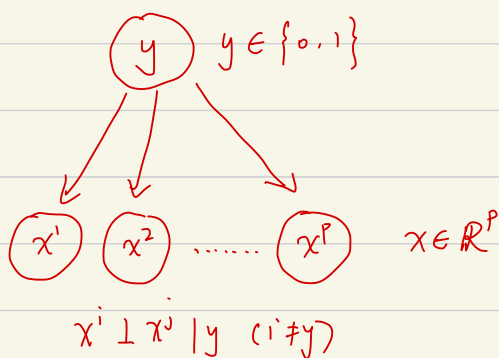
只不过我自己手推了 $x^T A x$ 对 A 的导数, 直接使其一步到位.

原视频中的感觉更巧妙.

朴素贝叶斯 (Naive Bayes)

思想: 朴素贝叶斯假设 (条件独立性假设) 就是为了简化运算
最简单的概率图, 也是最简单的概率生成模型
(有向图).

→ 相同类别中, 不同特征之间是独立的.



$$\text{即 } p(x|y) = \prod_{i=1}^P p(x^i|y)$$

则对于给定数据:

$$\text{Data: } \{(x_i, y_i)\}, x_i \in \mathbb{R}^P, y_i \in \{0, 1\}.$$

$$\begin{aligned} \text{MAP: } \hat{y} &= \underset{y}{\operatorname{argmax}} p(y|x) \\ &= \underset{y}{\operatorname{argmax}} \frac{p(x|y) p(y)}{p(x)} \\ &= \underset{y}{\operatorname{argmax}} p(x|y) p(y) \end{aligned}$$

对于 $p(y)$: 2 分类: $y \sim \text{Bernoulli Distribution}$

多分类: $y \sim \text{Categorical Distribution}$

对于 $p(x|y)$: $p(x|y) = \prod_{i=1}^P p(x_i|y)$.

当 x 离散: $x_i \sim \text{Categorical Distribution}$

当 x 连续: $x_i \sim N(\mu_i, \sigma_i^2)$

对上述提到的分布进行简单介绍:

Categorical Distribution:

这是将伯努利分布推广到多分类的结果: 意味着 P 个 x_i 中, 只有一个为 1, 其他为 0

设 X 有 P 个取值: (x_1, x_2, \dots, x_P) , $x_i \in \{0, 1\}$ 且 $\sum x_i = 1$.

每个取值对应的概率: $(\alpha_1, \alpha_2, \dots, \alpha_P)$, $\alpha_i \in [0, 1]$, $\sum \alpha_i = 1$

$$x_i | P(X = x_i) = \prod_{i=1}^P \alpha_i^{x_i}$$

Bernoulli Distribution $\xrightarrow{n \text{ 次}}$ Binomial Distribution

\downarrow 多分类

Categorical Distribution $\xrightarrow{n \text{ 次}}$ Multinomial Distribution.

下面将假设 $y \sim \text{Bernoulli Distribution}$

$x_i \sim \text{Categorical Distribution}$

对参数进行估计.

由于分布指定, 参数为确定之数, 则采用 MLE.

$$\hat{y} = \operatorname{argmax}_y p(x, y).$$

$$\mathcal{L}(\beta) = \log \prod_{i=1}^N P(x_i, y_i)$$

$$= \sum_{i=1}^N \log P(x_i | y_i) P(y_i)$$

$$= \sum_{i=1}^N \left[\log P(y_i) + \log \prod_{j=1}^P P(x_i^j | y_i) \right]$$

$$= \sum_{i=1}^N \left[\log P(y_i) + \sum_{j=1}^P \log P(x_i^j | y_i) \right].$$

由于 $y \sim \text{Bernoulli Distribution} \Rightarrow p(y) = \phi^y (1-\phi)^{1-y}$

$x^i \sim \text{Categorical Distribution} \Rightarrow p(x^i | y) = \prod_{k=1}^m \theta_k^{x_i^{j,k}}$

$\sum_{k=1}^m \theta_k = 1$
意思是 x 第 j 个特征取第 k 个值, $x_i^{j,k} \in \{0, 1\}$
0 代表取, 1 代表不取.

$$\text{则 } \mathcal{L}(\beta) = \sum_{i=1}^N \left[\log \phi^y (1-\phi)^{1-y} + \sum_{j=1}^P \log \prod_{k=1}^m \theta_k^{x_i^{j,k}} \right]$$

$$= \sum_{i=1}^N \left[\log \phi^{y_i} (1-\phi)^{1-y_i} + \sum_{j=1}^P \sum_{k=1}^m x_i^{j,k} \log \theta_k \right]$$

且不同特征的 θ_k 都不同.

令 $\beta = (\phi, \theta_k, \dots)$ 注意, 每个特征都有 k 个 θ_k , 但这里由于公式中的 θ_k 上标太多了, 且每个特征的 θ_k 求解方式都是一样的.
故没有给他区分

$$\beta = \operatorname{argmax}_{\beta} \mathcal{L}(\beta).$$

对于 ϕ .

$$\frac{\partial L(\beta)}{\partial \phi} = \frac{\partial}{\partial \phi} \sum_{i=1}^N \log \phi^{y_i} (1-\phi)^{1-y_i} = 0$$

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

对于 θ_k .

$$L(\beta) = \sum_{i=1}^N \sum_{j=1}^P \sum_{k=1}^m x_i^{j,m} \log \theta_k.$$

由于每个特征的 θ_k 都是不同的. 其实这一步要等到后面求偏导时可以把 $\sum_{j=1}^P$ 消去, 但是提前做也影响不了什么.

$$\text{则 } L(\beta) = \sum_{i=1}^N \sum_{k=1}^m x_i^m \log \theta_k.$$

因为本质上来说, 我当前只对该特征求解, 不关其他特征什么事.

由于 $x_i \in \{0, 1\}$. 设对于整个 Data, 在该特征中取第 k 个值的样本数为 N_k .

$$\text{则 } L(\beta) = \sum_{k=1}^m N_k \log \theta_k \quad \text{且 } \sum_{k=1}^m \theta_k = 1.$$

$$\text{用拉格朗日乘子法. } \lambda \quad F = \sum_{k=1}^m N_k \log \theta_k + \lambda (1 - \sum_{k=1}^m \theta_k).$$

$$\begin{cases} \frac{\partial F}{\partial \theta_k} = \frac{N_k}{\theta_k} - \lambda = 0 & \textcircled{1} \\ \frac{\partial F}{\partial \lambda} = 1 - \sum_{k=1}^m \theta_k = 0 & \textcircled{2} \end{cases}$$

由①. 得 $\theta_k = \frac{N_k}{\lambda}$ ③

将③代入②. 得: $\lambda = N$.

则 $\boxed{\theta_k = \frac{N_k}{N}}$ 对其他特征求法相同.