

Ch2. 统计判别

1. 通过对被识别对象进行多次观察及测量, 抽取特征构成特征向量, 并将其作为判别规则的输入。
对样本进行分类。

2. 在现实生活中, 只有通过大量观察, 其结果可能具有规律性。
3. 基于上述情况, 我们抽取出的特征向量实际是随机向量。

利用 Bayes 判别进行分类

欲确定样本 x 属于 w_1 类 or w_2 类, 则要看它属于 w_1 的概率大还是属于 w_2 的概率大。

即: if: $P(w_1|x) > P(w_2|x)$ 则 $x \in w_1$

反之, $x \in w_2$

而由 Bayes 公式 后验概率 $P(w_i|x) = \frac{P(x|w_i)P(w_i)}{P(x)}$?

$$= \frac{P(x|w_1)P(w_1)}{\sum_{i=1}^2 P(x|w_i)P(w_i)} \quad (2)$$

将 ② 代入 ①, 有:

若 $P(x|w_1) \cdot P(w_1) > P(x|w_2) \cdot P(w_2)$
则 $x \in w_1$; 反之, $x \in w_2$

再将 ③ 变形:

$$\text{若 } l_{12}(x) = \frac{P(x|w_1)}{P(x|w_2)} > \frac{P(w_2)}{P(w_1)} = \theta_{21}$$

↑ 似然阈值

则 $x \in w_1$; 反之 $x \in w_2$.

二. Bayes 最小风险判别
 Bayes 决策过程实际上是一个总风险的优化过程
 我们选择将条件风险最小化来使得预期损失最小化

① M 类分类问题的条件平均风险

$$r_j(x) = \sum_{i=1}^M L_{ij} P(W_i|x) \quad (*)$$

L_{ij} 是本来属于 W_i 类却判成了 W_j 类的逻辑代价

$$L_{ij} = \begin{cases} 0, & i=j \\ 1, & i \neq j \end{cases}$$

若对每一个模式 x 都计算出全部类别的平均风险值 $r_1(x), r_2(x), \dots, r_M(x)$, 并且将 x 指定为具有最小风险的那一类, 则这种分类器称之为最小平均风险分类器, 它属于:

$$r_j(x) = \frac{1}{p(x)} \sum_{i=1}^M L_{ij} p(x|W_i) \cdot p(W_i)$$

(由*经 Bayes 公式而来), 可简化为:

$$r_j(x) = \sum_{i=1}^M L_{ij} p(x|W_i) \cdot p(W_i)$$

这种分类器是按平均条件风险作为标准判别

三. Normal Distribution 下的 Bayes 分类器

① M 种模式下的判别函数

$$d_i(x) = p(x|W_i) \cdot p(W_i) = \ln p(x|W_i) + \ln p(W_i)$$

将概率密度代入

$$\ln p(W_i) - \frac{1}{2} \ln |C_i| - \frac{1}{2} (x - m_i)^T C_i^{-1} (x - m_i)$$

其中 $m_i = E_i(x)$ $C_i = E_i[(x - m_i)(x - m_i)^T]$ $i=1, 2, \dots, M$

$d(x)$ 可进一步写为 $d(x) = \vec{w}^T \vec{x}$

$\vec{x} = (x_1, x_2, \dots, x_n, 1)$ 叫增广模式向量
 $\vec{w} = (w_1, w_2, \dots, w_n, w_{n+1})$ 叫增广权向量

①

② 两类情况:

$$d(x) = \vec{w}^T \vec{x} \begin{cases} < 0, & x \in w_1 \\ > 0, & x \in w_2 \end{cases}$$

③ M 类情况 1: (w_i/w_j 两两分)

用 $d(x)$ 将 w_i 和 w_j 分开

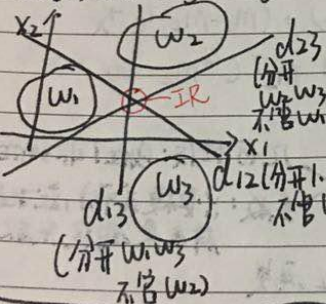
采用每对划分, 每个 $d(x)$

只能分开两个界面

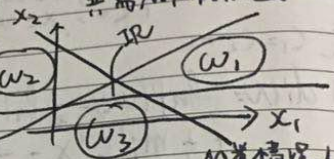
$$d_{ij}(x) = \vec{w}_{ij}^T \vec{x}$$

若 $d_{ij}(x) > 0, \forall i \neq j$, 则 $x \in w_i$

重要性: $d_{ij} = -d_{ji}$



④ M 类情况 2: (w_i/w_j 两两分)



若对某一区域, $d_{ij}(x) > 0$ 的条件超过 2 个, 则分类失败. 该区域称为不确定区域 (IR)

$$\text{共需 } C_M^2 = \frac{M(M-1)}{2} \text{ 个判别函数}$$

若 $\forall d_{ij}(x)$, 找不到 $\forall i \neq j, d_{ij}(x) > 0$ 则称为 IR.

③ M 类情况 3 (情况 2 的特例): 没有 IR 的情况

把 M 类问题分成 M-1 个两类问题 共需 M 个判别函数

把情况 2 中 $d_{ij}(x)$ 换成 $d_{ij}(x) - d_{ji}(x) = (w_i - w_j)^T x$ 则 $d_{ij}(x) > 0$

相当于 $d_{ij}(x) > d_{ji}(x), \forall i \neq j$

② 总结

若模式可用线性判别函数分开, 则称之为线性可分. 反之为线性不可分

GDA中我们假设 \$X\$ 服从多元正态分布。以 \$y=0\$ 为例，\$X\$ 服从 \$N(\mu_0, \Sigma)\$。同理 \$y=1\$ 时，\$X\$ 服从 \$N(\mu_1, \Sigma)\$。则有 \$p(X|y) = \frac{1}{(2\pi)^{d/2} |\Sigma|^{1/2}} \exp\left(-\frac{1}{2} (X - \mu_y)^T \Sigma^{-1} (X - \mu_y)\right)\$。

若直接对多元正态分布 \$p(X|y)\$ 建模，则有 \$2^d\$ 个不同值。则 \$p(X_1, \dots, X_{sw}|y) = p(X_1|y) \dots p(X_{sw}|y)\$。

在NB中我们假设 \$X\$ 的每一维都是独立的。则 \$\log\$ 似然为 \$\sum_{j=1}^n \log p(x_j|y) = \sum_{j=1}^n [y^{(j)} \log \phi_j + (1 - y^{(j)}) \log (1 - \phi_j)] + \sum_{j=1}^n [x_j^{(1)} \log \phi_j + (1 - x_j^{(1)}) \log (1 - \phi_j)]\$。

五. Bayes 模型 (得 \$\phi_j = \frac{\sum_{i=1}^n y^{(i)} x_j^{(i)}}{\sum_{i=1}^n x_j^{(i)}}\$ 即所有 \$y=0\$ 的样本中有单词 \$X_j\$ 的邮件数除以 \$y=0\$ 的样本总数)

background: 估计 (未知) 预测 (未知) 未知样本估计 (未知) 未知样本估计 (未知)

Bayes 决策分类: choose = $\begin{cases} 1, & \text{if } (P(C=1|X) > P(C=0|X)) \\ 0, & \text{else} \end{cases}$

Bayes 公式: $P(C|X) = \frac{P(X|C) \cdot P(C)}{P(X)} \propto P(X|C) P(C)$

Loss: λ_{ik} : 样本属于 \$k\$，却分到 \$i\$ 的风险损失
 α_i : 分到 \$i\$
 选取 \$\alpha_i\$ 的风险 $R(\alpha_i|X) = \sum_{k=1}^K \lambda_{ik} P(C_k|X)$
 Choose \$\alpha_i\$, if $R(\alpha_i|X) = \min_K R(\alpha_k|X)$

0-1 Loss: $\lambda_{ik} = \begin{cases} 0, & i=k \\ 1, & i \neq k \end{cases} = \sum_{k \neq i} P(C_k|X)$
 则 $R(\alpha_i|X) = \sum_{k \neq i} P(C_k|X) = 1 - P(C_i|X)$
 Choose \$\alpha_i\$ if $P(C_i|X) = \max P(C_i|X)$

分类的一个方法是用判别函数 \$g(x)\$
 Choose \$C_i\$ if $g_i(x) = \max_K g_K(x)$
 例, \$g_i(x)\$ 可取 $-\log \frac{P(C_i|X)}{P(X|C_i)P(C_i)}$

假设检验
No. of
p-value

中假说
y=0

y=1

y=0

y=1

y=0

y=1

y=0

y=1

y=0

y=1

y=0

y=1

y=0

y=1

y=0

y=1

y=0

y=1

y=0

y=1

y=0

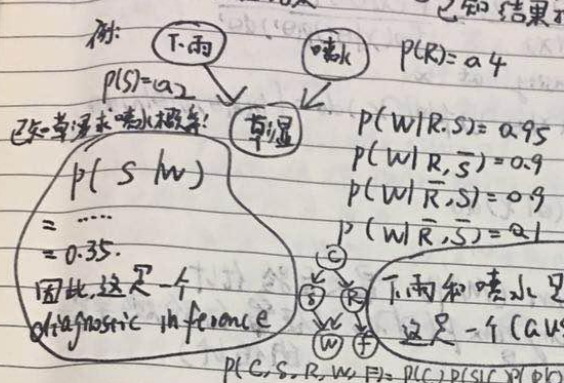
y=1

y=0

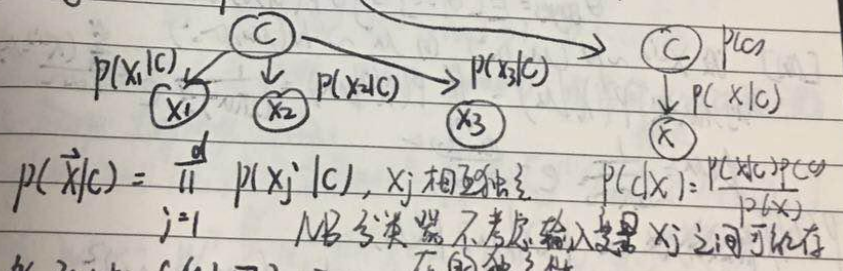
y=1

对于2类判别问题, 我们可令 $g(x) = g_1(x) - g_2(x)$
那么, $choose = \begin{cases} c_1, & \text{if } g(x) > 0 \\ c_2, & \text{else} \end{cases}$

Bayesian Networks 贝叶斯网: 其实就是一个图
Causal graph 因果图
Diagnostic inference 已知原因推结果
explaining away: 已知结果推原因



用贝叶斯网络分类 其实就是在 Diagnostic inference
Naive Bayes' classifier
- $P(x_1, \dots, x_d) = \prod_{i=1}^d P(x_i | \text{params}(x_i))$



最大似然估计 (MLE): 寻找 θ , 使得 $P(X|\theta)$ 中的样本 $X^{(i)}$ 最有可能出现
对于样本 $X = \{X^{(i)}\}_{i=1}^N$, 选择 θ 使得 $P(X|\theta)$ 最大, 即为 MLE

似然函数 $L(\theta|X) = P(X|\theta) = \prod_{i=1}^N P(X^{(i)}|\theta)$

$$\mathcal{L}(\theta|X) = \log L(\theta|X) = \sum_{i=1}^N \log p(X^{(i)}|\theta)$$

$$\hat{\theta} = \arg \max_{\theta} \mathcal{L}(\theta|X)$$

对于样本 X , 参数 θ 为多少时 X 出现的概率最大? 此时 θ 就是 MLE 估计的 θ 。

贝叶斯估计.
MLE 认为 θ 是固定值. 而 Bayes 认为 θ 有先验分布 $p(\theta)$

给定 $p(\theta|x) = \frac{p(x|\theta)p(\theta)}{p(x)} = \frac{p(x|\theta)p(\theta)}{\int p(x|\theta')p(\theta')d\theta'}$

Estimation of density at x
Full Bayes approach: $p(x|x) = \int p(x|\theta, x)p(\theta|x)dx = \int p(x|\theta)p(\theta|x)dx$

prediction: $y = \int f(x|\theta)p(\theta|x)dx$

积分不好算, 称为:

Maximum a posteriori (MAP): 最大后验估计

$\theta_{map} = \arg \max_{\theta} p(\theta|x)$ (使用后验概率最大的估计)

MLE: 最大似然估计:

$\theta_{ML} = \arg \max_{\theta} p(x|\theta)$

Bayes' 贝叶斯估计:

$\theta_{Bayes} = E[\theta|x] = \int \theta p(\theta|x) d\theta$

[例] 设 $X^{(i)} \sim N(\mu, \sigma^2)$ 而 $\mu \sim N(\mu_0, \sigma_0^2)$

则似然 $p(X|\mu) = \prod_{i=1}^N p(X^{(i)}|\mu) = \frac{1}{(2\pi\sigma^2)^{N/2}} e^{-\frac{1}{2\sigma^2} \sum_{i=1}^N (X^{(i)} - \mu)^2}$

$p(\mu) = \frac{1}{\sqrt{2\pi\sigma_0^2}} e^{-\frac{(\mu - \mu_0)^2}{2\sigma_0^2}}$

$p(\mu|x) \propto p(x|\mu)p(\mu)$

只考虑指数: $P(\mu|x) = -\frac{1}{2\sigma^2} (\mu - \mu_N)^2$

$= -\frac{1}{2\sigma^2} (\mu^2 - 2\mu_N\mu + \mu_N^2)$

$= -\frac{1}{2\sigma^2} \mu^2 + \frac{\mu_N}{\sigma^2} \mu + \text{Const}$

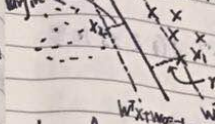
$= -\frac{\mu^2}{2} \left(\frac{1}{\sigma_0^2} + \frac{N}{\sigma^2} \right) + \mu \left(\frac{\mu_0}{\sigma_0^2} + \frac{1}{\sigma^2} \sum_{i=1}^N X^{(i)} \right) + \text{Const}$

$\therefore \mu_N = \frac{\sigma^2}{N\sigma_0^2 + \sigma^2} \mu_0 + \frac{N\sigma_0^2}{N\sigma_0^2 + \sigma^2} \mu_{ML} \rightarrow \frac{1}{N} \sum_{i=1}^N X^{(i)}$

$\frac{1}{\sigma_N^2} = \frac{1}{\sigma_0^2} + \frac{N}{\sigma^2}$

支持向量机 (SVM)

① 硬边界 SVM



定义 hyper plane:

再定义两条线

设这两个点分

别: $W^T X$

要最大化 margin

约束条件: W

转为优化:

这是一个二

维优化问题

$Lp(w_1, w_2)$

$= \frac{1}{2} w^T w$

$\begin{cases} \frac{\partial L}{\partial w} = 0 \\ \frac{\partial L}{\partial w_0} = 0 \end{cases}$

$\frac{\partial L}{\partial w_0} = 0$

$Ld(\{x_i\})$

$= -\frac{1}{2}$