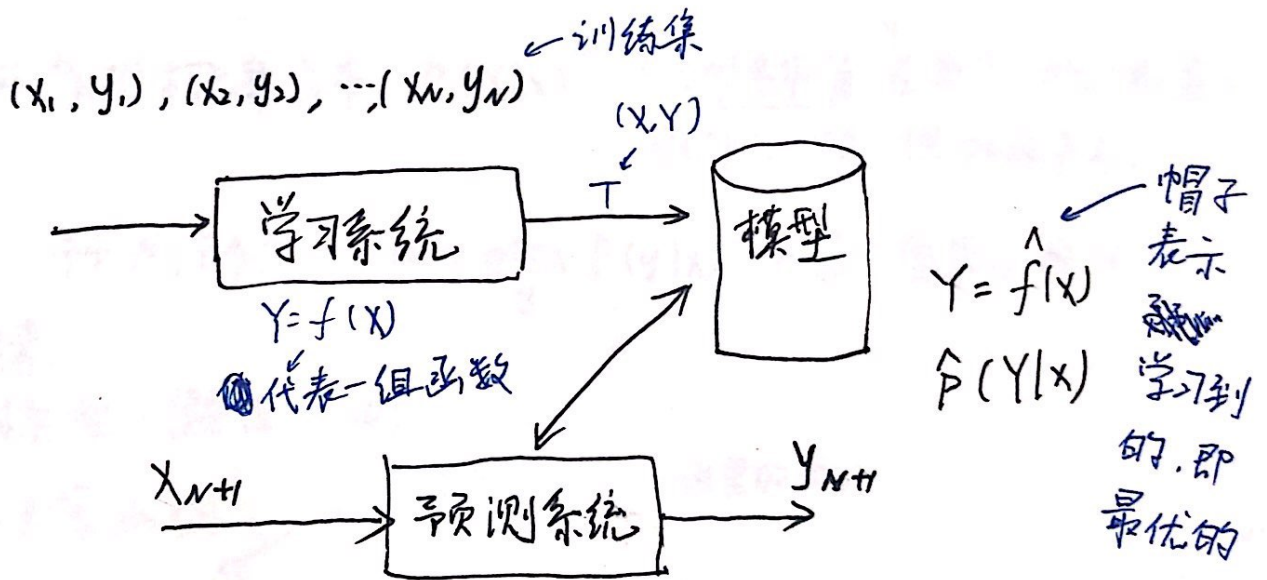


1.1 统计学习

监督学习的实现步骤

1. 有限训练集合
 2. 解设模型假设空间, 即备选模型确定
 3. 确定模型选择准则, 即学习策略
 4. 实现求解最优模型的算法
 5. 通过学习方法选择最优模型
 6. 利用学习的最优模型对新数据进行预测or分析
- 学习系统



1.2 监督学习

拿到的数据有随机性,

∴ 用联合概率分布表示

训练集

$$T = \{(x_1, y_1), (x_2, y_2), \dots, (x_n, y_n)\}$$

实例 x 的特征向量

$$x = (x^{(1)}, x^{(2)}, \dots, x^{(n)})^T$$

一般输入空间 = 特征空间

但, 若 x 为 输入空间, 但我们要将 $x \rightarrow x^2 + x^3$, 则

x^2 与 x^3 为 特征空间

模型:

1) 决策函数 $Y = f(x)$

表示一组假设空间的函数

预测形式 $y = \hat{f}(x)$

最优的, 就一个值

2) 条件概率分布 $P(Y|x)$ ^{表示} Y 有很多取值, 在每个不同的取值上概率不一样, 但加起来为 1

预测形式 $\arg \max_y P(y|x)$ 取最大值最应的 y .

三要素

① 模型 (假设空间)

决策函数

代表每个函数 (一组里的一个)

$$F = \{f \mid Y = f_\theta(x), \theta \in \mathbb{R}^n\}$$

给定假设空间, 求出 θ .

即求出了最优.

条件概率分布

$$F = \{P \mid P_\theta(Y|x), \theta \in \mathbb{R}^n\}$$

例: $Y = a_0 + a_1 X, \theta = (a_0, a_1)^T$

例: $Y \sim N(a_0 + a_1 X, \sigma^2),$

$\theta = (a_0, a_1)^T$

② 策略:

损失函数 (针对每个实例说的)

针对决策函数

$$L(Y, f(x)) = \begin{cases} 1, & Y \neq f(x) \\ 0, & Y = f(x) \end{cases} \quad (\text{分类问题})$$

真实值 预测值

$$L(Y, f(x)) = (Y - f(x))^2 \quad (\text{回归问题})$$

平方损失 (对差值大的实例敏感)

$$L(Y, f(x)) = |Y - f(x)| \quad (\text{回归问题})$$

绝对值损失

针对条件

概率分布

$$L(Y, P(Y|X)) = -\log P(Y|X)$$

对数似然函数

经验风险最小化:

$$\min_{f \in F} \frac{1}{N} \sum_{i=1}^N L(y_i, f(x_i))$$

结构风险最小化

$$\min_{f \in F} \frac{1}{N} \sum_{i=1}^N L(y_i, f(x_i)) + \lambda J(f)$$

正则项

衡量数据集上经验与模型复杂度重要程度

表示这个模型的复杂度

力求 经验风险小 + 模型复杂度低

1.4 模型评估与模型选择

训练误差 (在训练集上)

$$\frac{1}{N} \sum_{i=1}^N L(y_i, \hat{f}(x_i))$$

测试误差 (在测试集上)

$$\frac{1}{N'} \sum_{i=1}^{N'} L(y_i, \hat{f}(x_i))$$

过拟合

例: 多项式拟合问题 (书上 M 次多项式拟合)

M 次项增加会导致“噪声”增加, \therefore 力求模型复杂度低

1.5 正则化与交叉验证 (模型选择的2个方法)

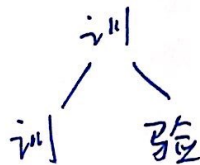
考虑

正则化:

最小化结构风险

$$\frac{1}{N} \sum_{i=1}^N L(y_i, f(x_i)) + \lambda J(f)$$

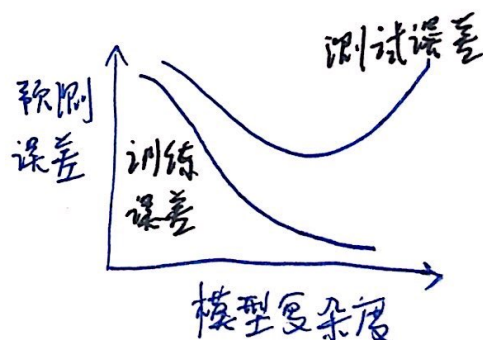
交叉验证



$M=1$

$M=3$

$M=9$



1.6 泛化能力

定理 1.1 泛化误差上界

对于二分类问题, 当假设空间是有限个函数的集合

$F = \{f_1, f_2, \dots, f_d\}$ 时, 对任意一个函数 $f \in F$, 至少以概率

$1 - \delta$, 以下不等式成立:

$$R(f) \leq \hat{R}(f) + \varepsilon(d, N, \delta)$$

念 yī pè xiū

期望风险 经验风险
(在所有数据上) (在训练数据集上)

其中

$$\varepsilon(d, N, \delta) = \sqrt{\frac{1}{2N} (\log d + \log \frac{1}{\delta})}$$

备选模
型的个
数

$d \uparrow, R(f) \uparrow$

样本量 $\uparrow, R(f) \downarrow$

1.7 生成模型与判别模型

生成方法 (X与Y都是随机变量)

$$P(Y|X) = \frac{P(X,Y)}{P(X)}$$

给定X, 不仅要学习Y的概率分布, 还要学习X与Y的联合概率分布.

判别方法 (不考虑X的随机)

$$f(X) \text{ 或 } P(Y|X)$$

(给定X, 学习 $f(X)$ 或Y的概率分布)

1.8 分类问题 (分类问题中的模型叫分类器)

TP — 将正类⁽¹⁾预测为正类 (2分类代表2个类型)

FN — 正⁽⁰⁾ → 负

FP — 负 → 正

TN — 负 → 负

精确率

$$P = \frac{TP}{TP+FP}$$

召回率

$$R = \frac{TP}{TP+FN}$$

1.9 标注问题

输入

$$X = (x^{(1)}, x^{(2)}, \dots, x^{(n)})^T$$

输出

$$y = (y^{(1)}, y^{(2)}, \dots, y^{(n)})^T$$

例：文本分类

输入：At Microsoft Research

$$X = \begin{pmatrix} \text{At} & \text{Microsoft} & \text{Research} \\ \parallel & \parallel & \parallel \\ x^{(1)} & x^{(2)} & x^{(3)} \end{pmatrix}$$

输出：At / 0 Microsoft / B Research / E

└───┬───┘
Microsoft · Research

$$y = (y^{(1)}, y^{(2)}, y^{(3)})$$

0 B E

意义

可出将文本分为一个一个词组

1.10 回归问题

输出为连续的值

★ 极大似然估计 ~~极大似然估计~~

例: 掷硬币, 设出现正面向上的概率为 θ

$$X_i = \begin{cases} 1, & \text{正} \\ 0, & \text{反} \end{cases}, X_i \sim b(1, \theta), P(X=x) = \theta^x (1-\theta)^{1-x}$$

∵ 每次掷硬币独立,

$$L(\theta) = P(X_1=x_1|\theta) \cdots P(X_n=x_n|\theta) \quad \therefore \text{可以把联合概率分布写成连乘形式}$$

$$= \prod_{i=1}^n \theta^{x_i} (1-\theta)^{1-x_i}$$

极大似然目的:

想找到 θ , 使得样本出现的概率最大, 即要最大化这个似然函数.

$$\begin{aligned} \max \ln L(\theta) &= \sum [\ln \theta^{x_i} + \ln (1-\theta)^{1-x_i}] \\ &= \sum x_i \ln \theta + (n - \sum x_i) \ln (1-\theta) \end{aligned}$$

求导

$$\frac{\partial \ln L(\theta)}{\partial \theta} = \frac{\sum x_i}{\theta} - \frac{n - \sum x_i}{1-\theta} = 0$$

等价于:

最大化似然函数的对数.

求得 $\boxed{\hat{\theta} = \frac{\sum x_i}{n}}$

步骤: ① 根据样本概率分布, 写出样本的联合概率似然

② 通过最大化似然函数求得参数估计值 函数

★ 极大似然估计 (完全根据样本信息)

★ 贝叶斯估计 (不仅根据样本信息, 还有先验信息)

∴ 一开始有了 $\pi(\theta)$

根据 x_1, \dots, x_n (即样本)

去调整对 θ 概率分布判断

即

$$P(\theta | x_1, \dots, x_n)$$

← 条件

← 联合密度

$$= \frac{P(\theta, x_1, \dots, x_n)}{P(x_1, \dots, x_n)}$$

← 边缘密度

(根据乘法公式)

$$= \frac{\pi(\theta) \cdot P(x_1 | \theta) \cdots P(x_n | \theta)}{\int P(\theta, x_1, \dots, x_n) d\theta}$$

← 联合密度对 θ 求积分

将 $\pi(\theta)$ 代入. 为了简化, 写成正比例形式 \propto (只写与 θ 有关的)

$$= \propto \theta^{\alpha-1} (1-\theta)^{\beta-1} \cdot \prod \theta^{x_i} (1-\theta)^{1-x_i}$$

$$= \theta^{\sum x_i + \alpha - 1}$$

$$(1-\theta)^{n - \sum x_i + \beta}$$

参数为

$\sum x_i + \alpha$, 与

$n - \sum x_i + \beta$

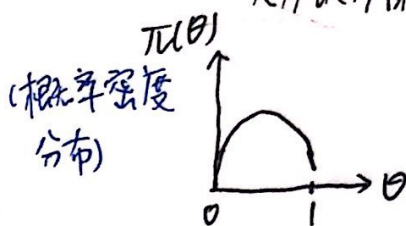
步骤:

① 根据参数, 给定后验信息的样本分布

② 给 θ 一个具体的值, 找一个使得后验分布最大的值, 即

找 β 分布的众数, 为 $\hat{\theta} = \frac{\sum x_i + \alpha - 1}{n + \alpha + \beta - 2}$

人为认为的他应该是这样的分布



一般 $[0, 1]$ 之间均匀概率分布可以用 β 分布表示.

伽马函数

$$\therefore \pi(\theta) = \frac{\Gamma(\alpha + \beta)}{\Gamma(\alpha) \Gamma(\beta)} \theta^{\alpha-1} (1-\theta)^{\beta-1}$$

∴ 分布由 α 与 β 决定

对比极大似然估计与贝叶斯估计

$$\text{MLE: } \hat{\theta} = \frac{\sum x_i}{n}$$

$$\text{Bayes: } \hat{\theta} = \frac{\sum x_i + \alpha - 1}{n + \alpha + \beta - 2} \quad \text{① 当 } n \rightarrow \infty \quad \text{样本个数.}$$

$$\text{Bayes} \rightarrow \frac{\sum x_i}{n}$$

$$\text{② 当 } n=1$$

$$\text{MLE} \rightarrow 1$$

$$\text{Bayes} \rightarrow \frac{\alpha}{\alpha + \beta - 1}$$

意义: 若样本量大,

先验分布微不足道.

若样本量小,

Bayes 不会出现极端情况

★ 泛化误差上界

对于二分类问题, 当假设空间是有限个函数的集合

$F = \{f_1, f_2, \dots, f_d\}$ 时, 对任意一个函数 $f \in F$, 至少以概率

$1 - \delta$, 以下不等式成立:

$$R(f) \leq \hat{R}(f) + \epsilon(d, N, \delta), \text{ 其中 } \epsilon(d, N, \delta) = \sqrt{\frac{1}{2N} (\log d + \log \frac{1}{\delta})}$$

↓
经验风险

Hoeffding 不等式

① 定义独立随机变量 x_1, x_2, \dots, x_n , $x_i \in [a_i, b_i]$

$$S_n = x_1 + x_2 + \dots + x_n, \quad E S_n = E(\sum x_i)$$

$$\textcircled{2} \quad P(S_n - E S_n \geq t) \leq \exp\left(-\frac{2t^2}{\sum (b_i - a_i)^2}\right)$$

损失函数取值

③ 不考虑随机变量的和, 考虑随机变量的均值

$$\bar{x}_n = \frac{S_n}{n}, \quad E(\bar{x}_n) = \frac{E S_n}{n}$$

可替换

$$\begin{aligned} P(\bar{x}_n - E(\bar{x}_n) \geq t) &= P(S_n - E S_n \geq nt) \\ P(E(\bar{x}) - \bar{x} \geq t) &\leq \exp\left(-\frac{2n^2 t^2}{\sum (b_i - a_i)^2}\right) \end{aligned}$$

↖ n^2 阶
↖ n 阶
↖ n 阶

$$\leq e^{-n} \leftarrow \bullet n \text{ 阶}$$

∴ 当样本量足够大, 概率为 0

④ $\hat{R}(f) = \frac{1}{N} \sum L(x_i, f(x_i)) \leftarrow \text{经验}$ $R(f) \leftarrow \text{期望}$

代入 Hoeffding 不等式

有 N 个候选函数

$$P(R(f) - \hat{R}(f) \geq t) \leq \exp\left(-\frac{2N^2 t^2}{N}\right) = \exp(-2N t^2)$$

(对于 1 个备选模型来说)

∴ $[a_i, b_i] = [0, 1] \therefore$ 为 N

等价于 不希望有一个备选模型, 使得期望风险到经验风险的距离大于 t . 即

$$\begin{aligned}
 & P(\exists f \in F: R(f) - \hat{R}(f) \geq t) \\
 &= P(R(f_1) - \hat{R}(f_1) \geq t \cup \dots \cup R(f_d) - \hat{R}(f_d) \geq t) \quad \leftarrow \text{并集} \\
 &\leq \sum_i P(R(f_i) - \hat{R}(f_i) \geq t) \\
 &\leq d \exp(-2Nt^2)
 \end{aligned}$$

等价于

$$P(\forall f \in F, R(f) - \hat{R}(f) \leq t) \geq 1 - \boxed{d \exp(-2Nt^2)}$$

\downarrow 设为 δ

解得 $t = \sqrt{\frac{1}{2N} (\log d + \log \frac{1}{\delta})}$ $\Sigma(d, N, \delta)$

\downarrow 书上为 ϵ