

统计学习方法

第1章 统计学习方法概论

1.1 统计学习

1.2 监督学习

1.3 统计学习三要素

1.3.1 模型

非概率模型:

$f = F^{(n)}$

概率模型:

$f = P(Y|X)$

1. 损失函数和风险函数

- (1) 0-1损失函数 (0-1 loss function)
- $$L(Y, f(X)) = \begin{cases} 1, & Y \neq f(X) \\ 0, & Y = f(X) \end{cases}$$
- (2) 平方损失函数 (quadratic loss function)
- $$L(Y, f(X)) = (Y - f(X))^2$$
- (3) 绝对损失函数 (absolute loss function)
- $$L(Y, f(X)) = |Y - f(X)|$$

(4) 对数损失函数 (logarithmic loss function) 或对数似然损失函数 (loglikelihood loss function)

$$L(Y, P(Y|X)) = -\log P(Y|X) \quad (1.8)$$

$R_{exp}(f) = E_L[L(Y, f(X))] = \int_{\mathcal{X} \times \mathcal{Y}} L(Y, f(X)) P(\mathbf{x}, y) d\mathbf{p} \quad (1.9)$
这是理论上模型f关于联合分布P(X,Y)的平均意义上的损失, 称为风险函数 (risk function) 或期望损失 (expected loss)。

模型f(X)关于训练数据集的平均损失称为经验风险 (empirical risk) 或经验损失 (empirical loss), 记作 R_{emp} 。

$$R_{emp}(f) = \frac{1}{N} \sum_{i=1}^N L(y_i, f(x_i)) \quad (1.10)$$

用训练数据集估计期望损失

2. 经验风险最小化与结构风险最小化

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

用, 比如, 极大似然估计 (maximum likelihood estimation) 就是经验风险最小化的一个例子。当模型是条件概率分布, 损失函数是对数损失函数时, 经验风险最小化或等价于极大似然估计。

但是, 当样本容量很小时, 经验风险最小化学习的效果就不太好, 会产生后面将要叙述的“过拟合(over-fitting)”现象。

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

结构风险最小化 (structural risk minimization, SRM) 是为了防止过拟合而提出来的策略。结构风险最小化等价于正则化 (regularization)。结构风险在经验风险上加上表示模型复杂度的正则化项 (regularization) 或惩罚项 (penalty term)。在假设空间、损失函数以及训练数据集确定的情况下, 结构风险的定义是

1.3.3 算法

1.4.1 训练误差与测试误差

1.4.2 过拟合与模型选择

1.4 模型评估与模型选择

1.5 正则化与交叉验证

1.5.1 正则化

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

以是参数向量的L₂范数:

$$L(w) = \frac{1}{N} \sum_{i=1}^N (f(x_i; w) - y_i)^2 + \frac{\lambda}{2} \|w\|^2$$

$f(x) = w^T \phi(x)$ 和 $\phi(x) = [1, x_1, \dots, x_n]^T$

1.5.2 交叉验证

1. 简单交叉验证
2. S折交叉验证
3. 留一交叉验证

1.6 泛化能力

1.6.1 泛化误差

首先给出泛化误差的定义。如果学到的模型是 \hat{f} , 那么用这个模型对未知数据预测的误差即为泛化误差 (generalization error)

$$R_{exp}(\hat{f}) = E_L[L(Y, \hat{f}(X))] = \int_{\mathcal{X} \times \mathcal{Y}} L(Y, \hat{f}(X)) P(\mathbf{x}, y) d\mathbf{p} \quad (1.20)$$

泛化误差反映了学习方法的泛化能力。如果一种学习方法学习的模型比另一种方法学习的模型具有更小的泛化误差, 那么这种方法就更好。事实上, 泛化误差就是所学习到的模型的风险。

1.6.2 泛化误差上界

定理 1.1 (泛化误差上界) 对二分类问题, 当假设空间是有有限个函数的集合 $\mathcal{F} = \{f_1, f_2, \dots, f_M\}$ 时, 对任意一个训练集 S , 至少以概率 $1 - \delta$, 以下不等式成立:

$$R(f) \leq \hat{R}(f) + \epsilon(d, N, \delta) \quad (1.25)$$

其中,

$$\epsilon(d, N, \delta) = \sqrt{\frac{1}{2N} (\log d + \log \frac{1}{\delta})} \quad (1.26)$$

1.7 生成模型与判别模型

监督学习方法又可以分为生成方法 (generative approach) 和判别方法 (discriminative approach)。所学到的模型分别称为生成模型 (generative model) 和判别模型 (discriminative model)。

生成方法由数据学习联合概率分布 $P(X, Y)$, 然后求出条件概率分布 $P(Y|X)$ 作为预测的模型, 即生成模型:

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

这样的方法之所以称为生成方法, 是因为模型表示了给定输入 X 产生输出 Y 的生成关系。典型的生成模型有: 朴素贝叶斯法和隐马尔可夫模型; 将在后面章节进行相关讨论。

判别方法由数据直接学习决策函数 $f(X)$ 或者条件概率分布 $P(Y|X)$ 作为预测的模型, 即判别模型。判别方法关心的是对给定的输入 X , 应该预测什么样的输出 Y 。典型的判别模型包括: k近邻法、感知机、决策树、逻辑斯谛回归模型、最大熵模型。支持向量机、 boosting 方法、神经网络等。将在后面章节讨论。

生成方法的特点: 生成方法可以还原出联合概率分布 $P(X, Y)$, 而判别方法则不能; 生成方法的学习收敛速度更快, 即当样本容量增加的时候, 学到的模型可以更有效地接近真实模型; 当存在未知变量时, 仍可以生成学习方法学习, 此时判别方法不能使用。

判别方法的特点: 判别方法直接学习的是条件概率 $P(Y|X)$ 或决策函数 $f(X)$, 直接面对预测, 因此学习的性能要好些; 由于直接学习 $P(Y|X)$ 或 $f(X)$, 可以对数据进行各种程度上的抽象、定义特征并使用特征, 因此可以简化学习问题。

1.8 分类问题

对于二分类问题常用的评价指标是精确率 (precision) 与召回率 (recall)。通常以关注的类为正类, 其他类为负类。分类器在测试数据集上的预测或正确或不正确, 4种情况出现的总分类别记作:

TP——将正类预测为正类数;
FN——将正类预测为负类数;
FP——将负类预测为正类数;
TN——将负类预测为负类数。

精确率定义为

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

召回率定义为

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

此外, 还有 F 值, 是精确率和召回率的调和均值, 即

$$F_1 = \frac{2}{\frac{1}{P} + \frac{1}{R}}$$

$$F_1 = \frac{2TP}{2TP + FP + FN}$$

许多统计学习方法可以用于分类, 包括k近邻法、感知机、朴素贝叶斯法、决策树、决策列表、逻辑斯谛回归模型、支持向量机、提升方法、贝叶斯网络、神经网络、Winnow等。本书将讲述其中一些主要方法。

1.9 标注问题

标注问题分为学习和标注两个过程 (如图 1.3 所示)。首先给定一个训练数据集 $T = \{(\mathbf{x}_1, y_1), (\mathbf{x}_2, y_2), \dots, (\mathbf{x}_N, y_N)\}$ 这里, $\mathbf{x}_i = (x_1^{(i)}, x_2^{(i)}, \dots, x_n^{(i)})^T$, $i=1, 2, \dots, N$, 是输入观测序列, $y_i = (y_1^{(i)}, y_2^{(i)}, \dots, y_m^{(i)})^T$ 是相应的输出标记序列, n 是序列的长度, 对不同样本可以有不同值。学习系统基于训练数据集构建一个模型, 表示为条件概率分布: $P(Y_1, Y_2, \dots, Y_m | X_1, X_2, \dots, X_n)$ 这里, 每一个 $X^{(i)}$ ($i=1, 2, \dots, n$) 取值为所有可能的观测, 每一个 $Y^{(i)}$ ($i=1, 2, \dots, m$) 取值为所有可能的标记, $-a \leq a \leq N$ 。标注系统按学习得到的条件概率分布模型, 对新的输入观测序列得到相应的输出标记序列。具体地, 对一给定的观测序列 $\mathbf{x} = (x_1^{(1)}, x_2^{(1)}, \dots, x_n^{(1)})^T$ 找到使条件概率 $P(Y_1^{(1)}, Y_2^{(1)}, \dots, Y_m^{(1)} | (x_1^{(1)}, x_2^{(1)}, \dots, x_n^{(1)})^T)$ 最大的标记序列 $\mathbf{y}_{opt} = (y_1^{(1)}, y_2^{(1)}, \dots, y_m^{(1)})^T$ 。

标注常用的统计学习方法有: 隐马尔可夫模型、条件随机场。

1.10 回归问题

习题

1.1 说明伯努利模型的最大似然估计以及贝叶斯估计中的统计学习方法三要素。伯努利模型是定义在取值 {0, 1} 的随机变量上的概率分布。假设观测到由伯努利模型生成的数据集生成数据, 求 θ 值的最大似然估计, 这时可以用极大似然估计或贝叶斯估计求出结果 $\hat{\theta}$ 的概率。

1.2 通过经验风险最小化推导极大似然估计, 证明模型是条件概率分布, 当损失函数是对数损失函数时, 经验风险最小化等价于极大似然估计。