

Chap01. 统计方法概论

I. 统计学习要素

$$X \rightarrow Y$$

↓

$$P(Y|X) \Rightarrow \text{分类}$$

$$\text{OR } Y = f(X)$$

↓ 回归分析

机器学习 =

$$\Rightarrow \text{模型} + \text{策略} + \text{算法}$$

1. 模型 \Rightarrow hypothesis space \Rightarrow 一般是有限的

$$\textcircled{1} \mathcal{F} = \{f | Y = f(X)\} \rightarrow \mathcal{F} = \{f | Y = f_0(X), 0 \in \mathbb{R}\}$$

$$\textcircled{2} \mathcal{P} = \{P | P(Y|X)\} \rightarrow \mathcal{F} = \{P | P_0(Y|X), 0 \in \mathbb{R}^n\}$$

①: 决策函数, 表示非概率模型

② 条件概率, 表示概率模型

2. 策略 \Rightarrow 怎么在假设空间中选取最优模型

(1). 损失函数和风险函数 [模型的好坏]

Loss function & cost function

$$\textcircled{1} \text{ 0-1 loss function: } L(Y, f(X)) = \begin{cases} 1, & Y \neq f(X) \\ 0, & Y = f(X) \end{cases} \quad (LF)$$

② Quadratic LF

② Quadratic Loss: $L(Y, f(x)) = (Y - f(x))^2$

③ absolute LF: $L(Y, f(x)) = |Y - f(x)|$

④ log LF: $L(Y, f(x)) = -\log p(Y|x)$

条件概率模型

Expected Loss =

$$R_{exp} = \iint_{x \times Y} \underbrace{p(x, Y)}_{\text{未知}} L(Y, f(x)) dx dy$$

Empirical Loss =

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

$\updownarrow N \rightarrow \infty$

(2). Empirical Risk Minimization (ERM)

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

N 较大时 over-fitting \rightarrow β 过

Structural Risk Minimization (SRM) Regularizer

\hookleftarrow

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

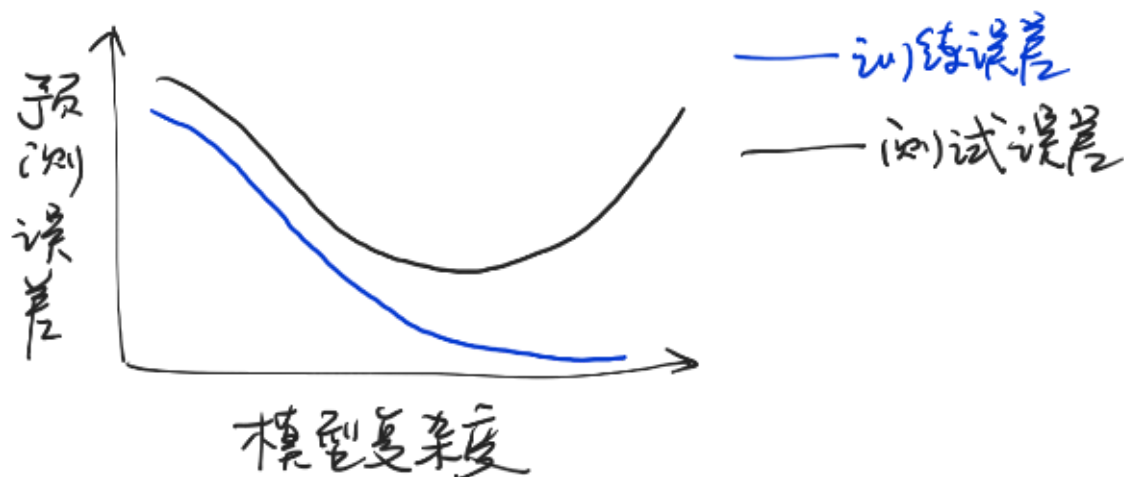
$J(f)$ 为模型的复杂度: f 越复杂, $J(f)$ 越大
如何表征复杂度?

$$f(x) \Rightarrow \min_{f \in \mathcal{F}} R_{SRM}(f)$$

3. 算法: 学习模型的具体方法, [高效求解与最优解]

II. 模型选择与正则化, 交叉验证

Over-fitting: 所选模型复杂度(参数个数)比“真”模型高, 导致, 虽然测试误差低, 但是预测误差大.

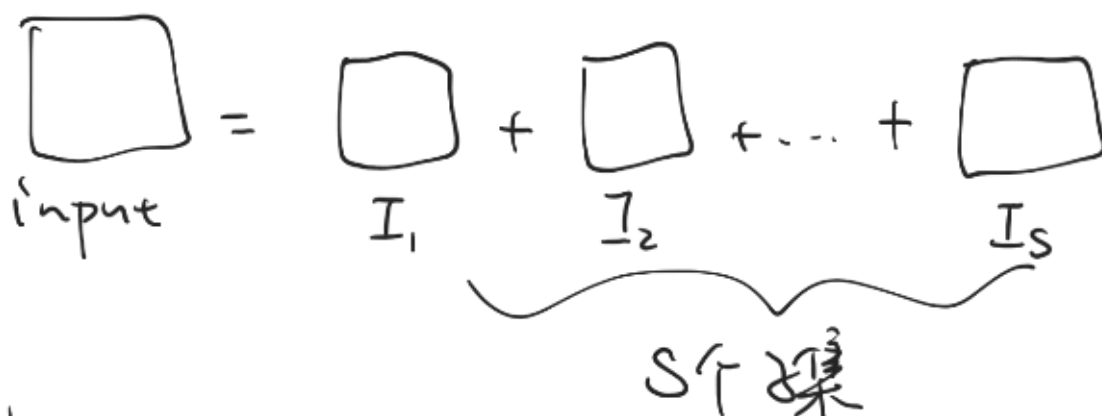


交叉验证: cross-validation.

① Simple cross-validation.



② S-fold cross-validation (应用较多)



每个子集 I_i 都包含 n/S 个样本, 且每个子集都包含 n/S 个测试样本.

1. 将 x_1, x_2, \dots, x_S 作为 test set, $x_i \neq x_j$
作为 training set. 每次的平均损失为 $L_i[f]$

(2) 最佳所选模型 $f \Rightarrow \min_{f \in F} \frac{1}{S} \sum_{i=1}^S L_i[f]$

(3) leave-one-out cross-validation ($S=1V$)

IV. 生成模型与判别模型

1. 生成模型: (generative model)

数据学习联合概率分布 $P(X, Y)$, 然后求出生成概率分布 $P(Y|X)$ 作为预测的模型.

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

例如: 朴素贝叶斯法和隐马尔可夫模型

2. 判别模型: (discriminative model)

数据直接学习决策函数 $f(x)$ 或条件概率分布 $P(Y|X)$ 作为预测的模型.

例如: k-邻近算法, 感知机, 决策树, 逻辑斯蒂回归模型, 最大熵模型, 支持向量机, 提升方法, 和条件随机场.

IV. 监督学习问题

1. 分类问题 = $\left\{ \begin{array}{l} Y \text{ 为有限个离散值} \\ X \text{ 可为连续或离散} \\ f(x) \parallel P(Y|x) \rightarrow \text{分类器} \\ 0-1 \text{ loss function} \end{array} \right.$

k-近邻法, 感知机, 朴素贝叶斯法, 决策树, 决策表
logistic 回归模型, 支持向量机, 提升方法, 贝叶斯网络,
神经网络, Winnow, ...

2. 标记问题 $\left\{ \begin{array}{l} Y \text{ 为标记序列} \\ X \text{ 为观测序列} \\ p(y^{(1)}, y^{(2)}, \dots, y^{(n)} | x^{(1)}, x^{(2)}, \dots, x^{(n)}) \end{array} \right.$

隐马尔可夫模型, 条件随机场

3. 回归问题 (含) 函数拟合 $\Rightarrow Y = f(X)$

Loss function: 常用为 MEAN-SQUARED-LOSS
可用最小二乘法求解, $f(x)$