

No.

Date.

chap 机器学习概述

模型

学习准则

优化算法

机器学习算法要素  $\Rightarrow$  统计推断问题

频率学派  $\Rightarrow$  固定参数

贝叶斯  $\Rightarrow$  随机变量, 且存在先验分布

基础

## 2.1 基本概念

特征向量  $\vec{x} = [x_1, x_2, \dots, x_n]^T$

训练集  $D = \{(\vec{x}^{(1)}, y^{(1)}), \dots, (\vec{x}^{(n)}, y^{(n)})\}$

标签  $y$

决策函数预测标签的值  $\hat{y} = f(\vec{x}, \theta)$

标签的条件概率  $p(y|\vec{x}) = p_y(\vec{x}, \theta)$

学习算法 A 找到一组参数  $\theta^*$ , 使得  $f(\vec{x}, \theta^*)$  近似真实关系

模型  $f(\vec{x}, \theta)$

### 2.1.1 指示函数

## 2.2 机器学习的基本要素

### 2.2.1 模型

$X$  输入空间

$(\vec{x}, y) \in X \times Y$

$Y$  输出空间

$p_y(y|\vec{x})$  真实条件概率分布

$F$  假设空间, 通常为参数化的函数族



No.

Date.

$\Rightarrow \Rightarrow \Rightarrow$  2.2.1 线性模型

$$f(\vec{x}, \theta) = \vec{w}^T \vec{x} + b$$

$\Rightarrow \Rightarrow \Rightarrow$  2.2.1.2 非线性模型

可以看成多个非线性基函数  $\phi(\vec{x})$  的线性组合.

$$f(\vec{x}, \theta) = \vec{w}^T \phi(\vec{x}) + b$$

$$\phi(\vec{x}) = [\phi_1(\vec{x}), \dots, \phi_k(\vec{x})]^T$$

$\Rightarrow \Rightarrow$  2.2.2 学习准则

$$|f(\vec{x}, \theta^*) - y| < \epsilon$$

$$|f_y(\vec{x}, \theta^*) - p_y(y|\vec{x})| < \epsilon$$

$R(\theta)$  期望风险

$$R(\theta) = E_{(\vec{x}, y) \sim P(\vec{x}, y)} [L(y, f(\vec{x}, \theta))]$$

$P(\vec{x}, y)$  真实的数据分布

$L(y, f(\vec{x}, \theta))$  损失函数

## 2.2.2.1 损失函数

## 0-1 损失函数

$$L(y, f(\vec{x}, \theta)) = \begin{cases} 0 & \text{if } y = f(\vec{x}, \theta) \\ 1 & \text{if } y \neq f(\vec{x}, \theta) \end{cases} = \mathbb{I}(y \neq f(\vec{x}, \theta))$$

## 平方损失函数

$$L(y, f(\vec{x}, \theta)) = \frac{1}{2} (y - f(\vec{x}, \theta))^2$$

## 交叉熵损失函数

$$L(y, f(\vec{x}, \theta)) =$$

假设  $y \in \{1, \dots, C\}$  个类别,  $f(\vec{x}, \theta) \in [0, 1]^C$

$p(y=c | \vec{x}, \theta) = f_c(\vec{x}, \theta)$  转化为标量的条件概率分布

## 交叉熵

$$L(y, f(\vec{x}, \theta)) = - \sum_{c=1}^C y_c \log f_c(\vec{x}, \theta)$$

例如:  $\vec{y} = [0, 0, 1]^T$ ,  $f(\vec{x}, \theta) = [0.3, 0.3, 0.4]^T$

$$L(\theta) = -(0 \times \log 0.3 + 0 \times \log 0.3 + 1 \times \log 0.4)$$

$$= -\log 0.4$$

## 对 one-hot 向量

$$L(y, f(\vec{x}, \theta)) = -\log f_y(\vec{x}, \theta)$$



No.

Date.

且  $f_{\theta}(\vec{x}, \theta)$  可以看作  $y$  的似然函数

负对数似然损失函数

Hinge 损失函数

$$\text{两分类问题} \begin{cases} -1 \\ +1 \end{cases}$$

$$L(y, f(\vec{x}, \theta)) = \max(0, 1 - yf(\vec{x}, \theta))$$

$$= [1 - yf(\vec{x}, \theta)]_+$$

$\Rightarrow \Rightarrow \Rightarrow$  2.2.2.2 风险最小化准则

$$R_D^{\text{emp}}(\theta) = \frac{1}{N} \sum_{i=1}^N L(y_i^{(w)}, f(x_i^{(w)}, \theta))$$

$$\theta^* = \arg \min_{\theta} R_D^{\text{emp}}(\theta)$$

ERM 经验风险最小化

过拟合

SRM 结构风险最小化

$$\theta^* = \arg \min_{\theta} R_D^{\text{struct}}(\theta)$$

$$= \arg \min_{\theta} R_0^{\text{emp}}(\theta) + \frac{1}{2} \lambda \|\theta\|^2$$

$$= \arg \min_{\theta} \frac{1}{N} \sum_{n=1}^N L(y^{(n)}, f(x^{(n)}, \theta)) + \frac{1}{2} \lambda \|\theta\|^2$$

 $\|\theta\|$  是  $L_2$  范数的正则化项 $\lambda$  正则化强度 $L_1$  正则化  $\lambda$  稀疏性

欠拟合

 $\Rightarrow$  2.2.3 优化算法

参数

超参数:

 $\Rightarrow \Rightarrow$  2.2.3.1 梯度下降法

$$\theta_{t+1} = \theta_t - \alpha \frac{\partial R_0(\theta)}{\partial \theta}$$

$$= \theta_t - \alpha \cdot \frac{1}{N} \sum_{n=1}^N \frac{\partial L(y^{(n)}, f(x^{(n)}, \theta))}{\partial \theta}$$

 $\alpha$  学习率



No.

Date.

⇒ ⇒ ⇒ 2.2.3.2 提前停止

验证集 从训练集中划分出来

验证集上错误率不再下降

⇒ ⇒ ⇒ 2.2.3.3 随机梯度下降法

SGD

$$\theta \leftarrow \theta - \alpha \frac{\partial L(\theta; (x^{(m)}, y^{(m)}))}{\partial \theta}$$

小批量梯度下降 Mini-Batch Gradient Descent

$$\theta_{t+1} \leftarrow \theta_t - \alpha \cdot \frac{1}{k} \sum \frac{\partial L(y, f(\vec{x}, \theta))}{\partial \theta}$$

## 2.3 机器学习的简单示例：线性回归

LR 自变量数量为1时称为简单回归

自变量...大于1时... 多元回归

模型  $f(\vec{x}; \vec{w}, b) = \vec{w}^T \vec{x} + b$

可写为  $f(\vec{x}; \vec{\tilde{w}}) = \vec{\tilde{w}}^T \vec{\tilde{x}}$

$\vec{\tilde{w}} = \vec{w} \oplus b = \begin{bmatrix} \vec{w} \\ b \end{bmatrix}$  增广权重向量

$\vec{\tilde{x}} = \vec{x} \oplus 1 = \begin{bmatrix} \vec{x} \\ 1 \end{bmatrix}$  增广特征向量

### ⇒ 2.3.1 参数学习

4种参数估计方法

ERM

SRM

最大似然估计

最大后验估计



⇒ 2.3.1.1 经验风险最小化

$$ER: R(w) = \sum_{n=1}^N L(y^{(n)}, f(\bar{x}^{(n)}, \bar{w}))$$

$$= \frac{1}{2} \sum_{n=1}^N (y^{(n)} - \bar{w}^T \bar{x}^{(n)})^2$$

$$= \frac{1}{2} \|\bar{y} - X^T \bar{w}\|^2$$

这里是个转置。

在这里

$$\bar{y} = \begin{bmatrix} y^{(1)} \\ \vdots \\ y^{(N)} \end{bmatrix}$$

$$\bar{w} = \begin{bmatrix} w_1 \\ 1 \\ w_N \\ b \end{bmatrix}$$

$$X = [\bar{x}^{(1)}, \dots, \bar{x}^{(N)}] =$$

$$\begin{bmatrix} x_1^{(1)} \\ \vdots \\ x_d^{(1)} \\ 1 \end{bmatrix}$$

$$x_1^{(N)}$$

$$x_d^{(N)}$$

$$1$$

如

$$\frac{\partial R(w)}{\partial \bar{w}} = \frac{1}{2} \frac{\partial \|\bar{y} - X^T \bar{w}\|^2}{\partial \bar{w}}$$



$$= \frac{1}{2} \cdot 2 (\bar{y} - X^T \bar{w}) \cdot \frac{\partial (\bar{y} - X^T \bar{w})}{\partial \bar{w}}$$

$$= (\bar{y} - X^T \bar{w}) \cdot (-X)$$

存疑

$$= -X (\bar{y} - X^T \bar{w})$$

$$\frac{1}{2} \frac{\partial R(\vec{w})}{\partial \vec{w}} = 0$$

$$-X(\vec{y} - X^T \vec{w}) = 0$$

$$X^T X (\vec{y} - X^T \vec{w}) = 0$$

$$X^T \vec{y} = X^T \vec{w}$$

$$X^T X \vec{y} = \vec{w} \quad \vec{w} = (X^T)^{-1} X^T \vec{y} = (X^T)^{-1} (X^T) X \vec{y}$$

LSE 最小二乘法估计

$$\vec{w} \leftarrow \vec{w} - 2X(\vec{y} - X^T \vec{w})$$

$\Rightarrow \Rightarrow$  2.3.1.2 SRM

岭回归 令  $XX^T$  对角线元素都加一个  $\lambda I$ , 使得  $(XX^T + \lambda I)$  不为 0

正则化

$$R(\vec{w}) = \frac{1}{2} \|\vec{y} - X^T \vec{w}\|^2 + \frac{1}{2} \lambda \|\vec{w}\|^2$$

$\Rightarrow \Rightarrow$  2.3.1.3 最大似然估计

假设  $y$  是一个随机变量服从均值为  $\vec{w}^T \vec{x}$  的正态分布  $\sigma^2$

$$p(y | \vec{x}, \vec{w}, \sigma) = N(y | \vec{w}^T \vec{x}, \sigma^2)$$

$$= \frac{1}{\sqrt{2\pi}\sigma} \exp\left(-\frac{(y - \vec{w}^T \vec{x})^2}{2\sigma^2}\right)$$



No.

Date.

似然函数

$$p(\vec{y} | X, \vec{w}, \sigma) = \prod_{n=1}^N p(y^{(n)} | \vec{x}^{(n)}, \vec{w}, \sigma)$$

$$= \prod_{n=1}^N \mathcal{N}(\vec{y}^{(n)} | \vec{w}^T \vec{x}^{(n)}, \sigma^2)$$

对数似然函数

$$\log p(\vec{y} | X, \vec{w}, \sigma) = \sum \log \mathcal{N}(\vec{y}^{(n)} | \vec{w}^T \vec{x}^{(n)}, \sigma^2)$$

MLE 最大似然估计

$$\text{令 } \frac{\partial \log p(\vec{y} | X, \vec{w}, \sigma)}{\partial \vec{w}} = 0 \Rightarrow \vec{w}^{ML} = (X^T X)^{-1} X^T \vec{y}$$

推导

$\Rightarrow \Rightarrow$  2.3.1.4 最大后验估计

假设  $\vec{w}$  为一个随机向量.

$p(\vec{w} | \vec{v})$  为各向同性的高斯分布

$$p(\vec{w} | \vec{v}) = \mathcal{N}(\vec{w} | 0, v^2 \mathbf{I}) \quad \text{--- } p(\vec{w})$$

$v^2$  为每一维上的方差

$\bar{w}$  的后验概率分布为

$$p(\bar{w} | x, \bar{y}, v, \sigma) = \frac{p(\bar{w}, \bar{y} | x, v, \sigma)}{\sum_{\bar{w}} p(\bar{w}, \bar{y} | x, v, \sigma)}$$

$$\propto p(\bar{y} | x, \bar{w}, \sigma) p(\bar{w} | v)$$

$\uparrow$   
 $\bar{w}$  的似然函数       $\bar{w}$  的先验

MAP

$$\bar{w}^{\text{MAP}} = \arg \max_{\bar{w}} p(\bar{y} | x, \bar{w}, \sigma) p(\bar{w} | v)$$

$$\log p(\bar{w} | x, \bar{y}, v, \sigma) \propto \log p(\bar{y} | x, \bar{w}, \sigma) + \log p(\bar{w} | v)$$

$\propto$

$\Rightarrow$



## 2.4 偏差与方差分解

$p_r(\vec{x}, y)$  真实分布  $f(\vec{x})$  模型

$$R(f) = E_{(\vec{x}, y) \sim p_r(\vec{x}, y)} [(y - f(\vec{x}))^2] \quad \text{期望误差}$$

$$f^*(\vec{x}) = E_{y \sim p_r(y|\vec{x})} [y] \quad \text{最优模型}$$

$$e = E_{(\vec{x}, y) \sim p_r(\vec{x}, y)} [(y - f^*(\vec{x}))^2] \quad \text{模型损失}$$

 $f(\vec{x})$  期望误差分解

$$R(f) = E_{(\vec{x}, y) \sim p_r(\vec{x}, y)} [(y - f(\vec{x}) + f^*(\vec{x}) - f^*(\vec{x}))^2]$$

$$= E_{\vec{x} \sim p_r(\vec{x})} [(f(\vec{x}) - f^*(\vec{x}))^2] + e \quad (2.63)$$

当前模型与真实模型之差距

$D$  上的期望差距  $f_0(\vec{x})$  与  $f^*(\vec{x})$

$$E_0[(f_0(\vec{x}) - f^*(\vec{x}))^2]$$

$$= E_0[(f_0(\vec{x}) - E_0[f_0(\vec{x})] + E_0[f_0(\vec{x})] - f^*(\vec{x}))^2] \quad (2.64)$$

$$= \underbrace{(E_0[E f_0(\vec{x})] - f^*(\vec{x}))^2}_{\text{偏差}^2} + E_0[(f_0(\vec{x}) - E_0[f_0(\vec{x})])^2] \quad (2.65)$$

方差

期望/错误可以分解为

$$R(f) = (\text{bias})^2 + \text{variance} + \epsilon$$

$$(\text{bias})^2 = E_{\vec{x}} \left[ (E[f_0(\vec{x})] - f^*(\vec{x}))^2 \right]$$

最优模型与模型均值 (12.67)

$$\text{variance} = E_{\vec{x}} \left[ E_0 \left[ (f_0(\vec{x}) - E_0[f_0(\vec{x})])^2 \right] \right]$$

模型与模型均值 (12.68)

当样本较多时, 方差较小, 我们可以选择能力强的模型减小偏差。

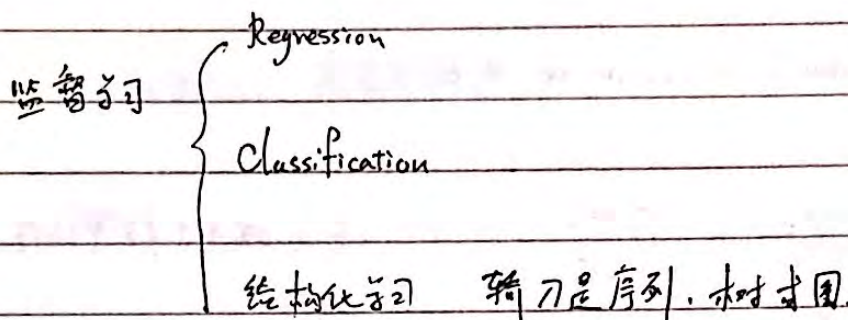
复杂模型。



No.

Date.

## 2.5 机器学习算法的类型



$\phi(\vec{x}, \vec{y})$  联合特征向量

结构化学习:

$$\hat{y} = \arg \max_{y \in G(\vec{x})} f(\phi(\vec{x}, y), \theta)$$

无监督学习 UL

聚类, 密度估计, 特征学习, 降维.

强化学习 RL

	监督学习	UL	RL
训练样本	$\{(\vec{x}^{(n)}, y^{(n)})\}$	$\{\vec{x}^{(n)}\}$	状态 $\vec{x}$ , 累积奖励 $G_t$
优化目标	$y = f(\vec{x})$ 或 $p(y \vec{x})$	$p(\vec{x})$ 或 $p(\vec{x} \vec{z})$	期望总回报 $E[G_t]$
训练准则	期望风险最小化 最大似然估计	最大似然估计 最小重构错误	策略评估 策略改进

## 2.6 数据的特征表示

### 图像特征

#### 文本特征

Bag-of-Words, Bow 模型

向量中的每一维  $x_i$  代表词表  $V$  中的第  $i$  个词是否在  $x$  中出现

每个连续词构成一个基本单元。二元特征 Bow

### 表示学习

⇒ 2.6.1 传统的特征学习

⇒ ⇒ 2.6.1.1 特征选择 FS

#### 子集搜索

贪心策略

前向搜索 添加该轮最优特征

反向搜索 删除最无用特征

过滤式

每次增加最有信息量的特征

信息增益

包裹式

每次增加对后续机器学习模型最有用特征

~~过滤式~~

$L_1$  正则化

稀疏特征



No.

Date.

## ⇒ ⇒ 2.6.1.2 特征抽取 FE

$$\bar{x}' = P \bar{x}$$

$\bar{x}$  原始特征向量  $\bar{x} \in R^d$

$\bar{x}'$  新特征向量  $\bar{x}' \in R^k$

$P$  映射矩阵  $P \in R^{k \times d}$

监督方法：抽取对一个特定的预测任务最有用特征

LDA

非监督方法：减少冗余信息和噪声。

PAC

	监督学习	无监督学习
FS	子集搜索 L1正则化 决策树	子集搜索
FE	LDA	PAC. 独立成分分析 流形学习. 自编码器

No.

Date. 2.7 评价指标

$T = (\bar{x}^{(1)}, y^{(1)}), (\bar{x}^{(2)}, y^{(2)}), \dots, (\bar{x}^{(N)}, y^{(N)})$  测试集

$\hat{Y} = \hat{y}^{(1)}, \dots, \hat{y}^{(N)}$  预测结果.

准确率

$$Acc = \frac{1}{N} \sum_{i=1}^N \mathbb{I}(y^{(i)} = \hat{y}^{(i)})$$

错误率

$$E = 1 - Acc$$

查准率和查全率

TP	FN	真的
FP	TN	假的
认为对的	认为错的	

查准率 Precision

$$P_c = \frac{TP_c}{TP_c + FP_c} \quad \frac{\text{真正对的}}{\text{所有认为对的}}$$

数据中真警察的比例

查全率 Recall

$$R_c = \frac{TP_c}{TP_c + FN_c} \quad \frac{\text{真对}}{\text{真类}}$$

~~数据中真警察的比例~~ 真警察在数据中的比例



F1值 是精确率和召回率的调和平均。

$$F_1 = \frac{2 \times P_c \times R_c}{P_c + R_c}$$

宏平均和微平均。

macro

micro

宏平均 每类性能指标的算术均值。

$$P_{macro} = \frac{1}{C} \sum P_c$$

$$R_{macro} =$$

$$F_{macro} =$$

AUC

ROC 曲线

PR 曲线

Top N

交叉验证

## 2.8 理论原理

## ⇒ 2.8.1 PAC 学习理论.

可能近似正确

可能

近似正确

同样条件下, 模型越复杂, 泛化误差越大.

换句话说, 模型越复杂越吃样本

## ⇒ 2.8.2 没有免费午餐定理 NFL

基于迭代的最优化算法, 不存在某种算法对所有问题都有效.

## ⇒ 2.8.3 引理.

## ⇒ 2.8.4 奥卡姆剃刀

最小描述长度.

$$\max_f \log p(f|D) = \max_f \log p(D|f) + \log p(f)$$

$$= \min_f -\log p(D|f) - \log p(f)$$

## ⇒ 2.8.5 归纳偏差 Inductive Bias

机器学习中的假设 贝叶斯中的先验 priors