

# 学习理论

偏差: Bias, 欠拟合, 高偏差

方差: variance, 过拟合, 高方差

泛化能力: 模型在测试集上的分类效果.

通过 联合界和 Hoeffding 不等式, 证明, 最优化 ERM (经验风险最小化) 能带来较小泛化误差

一致收敛性

$$E(h) \quad E(h^*)$$

一致收敛定理: 训练误差和泛化误差随样本数目的增大而更加接近

$$\text{最大差距为 } 2\gamma \quad E(h) \leq \hat{E}(h) + \gamma \leq \hat{E}(h^*) + \gamma \leq E(h^*) + \gamma + \gamma$$

$$\text{ERM 的定义: } \hat{E}(h) = \arg \min_{h \in H} \hat{E}(h)$$

训练误差集  
d. 的假设 h.

$$\hat{E}(h) = \frac{1}{m} \sum_{i=1}^m I(h(x^{(i)}) \neq y^{(i)})$$

训练误差

VC 维  $\rightarrow$  VC 维解释 SVM: SVM 算法会自动寻找一个具有较小 VC 维的假设类, 这样降低了 VC 维

ERM 的真定义:

模型选择:  $\begin{cases} \text{交叉检验} \begin{cases} \text{保留交叉检验} \\ \text{K重交叉检验} \end{cases} \\ \text{特征选择} \\ \text{防止过拟合} \end{cases}$  (通过减少模型参数, 简化模型来降低过拟合的发生)

正则化: 在文本分类问题上, 特征数有时会因为大于样本数, 这样就会产生过拟合. 添加正则化项是解决该问题的一个不错方法

在线学习:

对感知器算法,若正负样本分,那么在线学习算法也是收敛的.

当算法遇到瓶颈时,该选择什么样的方向来对算法进行改进?(算法诊断)

偏差/方差分析

一般的,高方差针对过拟合问题,即训练误差很小,但泛化误差很大.

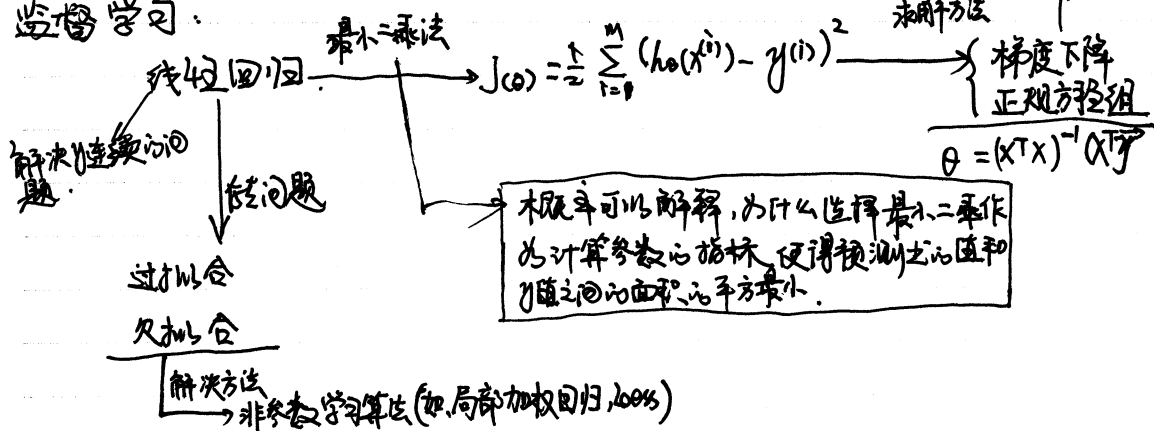
高偏差针对模型不合适问题,如特征数太少,表现是,训练误差和泛化误差都很大.

收敛与目标函数是否一致正确的判断

误差分析

收敛分析

监督学习:



分类问题算法: (线性分类)

Logistic回归 → 解决y是离散问题

可以使用以下方法拟合: 梯度上升, 牛顿方法 (更快, 但要计算Hessian矩阵)

伯努利分布 → logistic模型

多项式分布 → 线性模型

指数分布族: 若  $P(y|x; \theta)$ ,  $y$  为实数, 满足高斯分布, 得到基于最小二乘的线性回归

若  $y \in \{0, 1\}$ , 满足伯努利分布, 得到 logistic 回归

$$P(y; \eta) = \frac{b(\eta)}{Z(\eta)} \exp(\eta^T T(y) - a(\eta))$$

广义线性模型 (GLM)

假设: 1.  $y|x; \theta \sim \text{Exponential Family}(\eta)$ , 即假设讨论预测的变量  $y$ , 在给定  $x$  时, 自然参数的条件概率, 属于以  $\eta$  作为自然参数的指数分布族

2. 给定  $x$ , 目标是求出以  $x$  为条件的  $T(y)$  的期望  $E(T(y|x))$ , 即让学习算法输出:

$$h(x) = E(T(y|x))$$

3.  $\eta = \theta^T x$ , 即自然参数和输入特征之间的线性关系, 由  $\theta$  决定,  $T(y)$  是实数时才有意义.

若  $\eta$  是一个向量, 则  $\eta_i = \theta^T x_i$

给定一个指数分布族后, 推导 GLM: 求  $h(x)$ , 推导出  $\theta$

广义线性回归模型

判别学习算法：对  $P(y|x)$  建模

生成学习算法：对  $P(x|y)$  建模

→ 高斯判别 (GDA) → GDA 的假设很强  $\xrightarrow{\text{X}}$  Logistic 回归的假设较弱, 需要牺牲.

→ 朴素贝叶斯 → 事件问题 / 新样本通过 Laplace 平滑 解决

非线性分类算法

1. 神经网络

2. SVM (支持向量机) → 基础 → 最大间隔分类器.

↓ 为了更快地解决最优间隔问题

使用拉格朗日对偶性质

↓ 原始问题和对偶问题是等价解等价条件

KKT

↓ 得到最优间隔分类器的对偶形式

↓ 求解时发现目标函数中存在内积形式.

引入核函数 (得到完全全的 SVM 求解问题)

最优间隔分类器分: 线性: 一般不太符合实际

使用 非线性 软间隔分类器 方法, 得到最终问题

END.

使用序列最小化算法 (SMO) 求解

简单的与 SMO 同思想.

求出最优解的方法:  
坐标上升法

(之前介绍过梯度下降和平权法)

## 无监督学习

K-Means:

## 混合高斯分布 (MoG)

如: MoG中数据属于哪个分组可以看成2

一般地说, 任意形状的概率分布都可以用多个高斯分布函数去近似, 因而MoG应用广泛.

## 利用EM算法求解MoG

1. 设置初始值

2. E-Step, 估计隐含变量

3. M-Step, 根据2, 重新估计

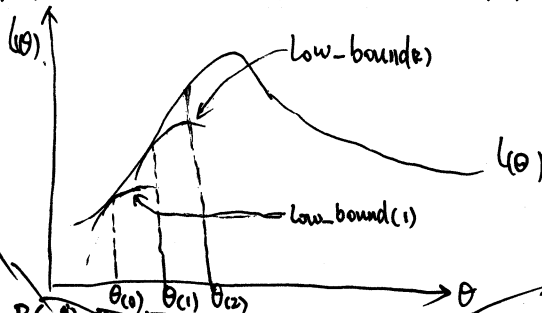
4. 重复2, 3, 直到参数变化  
小于阈值

→ 推导前提: Jensen不等式: 若  $f$  为凸函数, 则  $f'(x) \geq 0$ . 注意, 并不一定要求  $f$  为凸函数, 若存在二阶导数, 则  $f''(x) \geq 0$ . 再令  $x$  为随机变量, 则存在:  $f(E[x]) \leq E[f(x)]$

进一步, 若二阶导数恒大于0, 则不等式等号成立当且仅当  $x = E[x]$ , 即  $x$  是固定值.

若二阶导数<sup>不恒</sup>大于0, 则不等式的等号方向逆转.

思想: 存在一个不能直接求导的似然函数, 给定初始参数, 找到它初始参数下紧挨着似然函数的下界函数, 在下界函数上求极值来更新参数, 然后更新后的参数为初始值再次进行如上操作, 如图:



一般用于处理有隐含变量的模型.

此时目标函数为:

$$l(\theta) = \sum_{i=1}^m \log p(x^{(i)}; \theta) = \sum_{i=1}^m \log \sum_{z^{(i)}} p(x^{(i)}, z^{(i)}; \theta)$$

其实EM的一般化形式中, 可将目标函数看为:

$$J(\theta, Q) = \sum_{i=1}^m \sum_{z^{(i)}} Q(z^{(i)}) \log \frac{p(x^{(i)}, z^{(i)}; \theta)}{Q(z^{(i)})}$$

这样, EM算法就可以看做是对目标函数的坐标

上升过程: E-Step中,  $\theta$  不变, 调整  $Q$  使函数变大; 在 M-Step中,  $Q$  不变, 调整  $\theta$  使函数变大.

对  $l(\theta)$ , 对于当前  $\theta^{(t)}$ , 在下界上<sup>图</sup>进行极大似然估计得到  $\theta^{(t+1)}$ , 满足  $l(\theta^{(t+1)}) \geq l(\theta^{(t)})$ , 从而才能在下界函数上不断进行极大似然估计逼近真实值.

$$l(\theta^{(t+1)}) \geq \text{low\_bound}(\theta^{(t+1)}) \geq \text{low\_bound}(\theta^{(t)}) = l(\theta^{(t)})$$

假设条件

MoG: 混合高斯模型

MoNB: 混合贝叶斯模型

因子分解模型

在 EM 推导.

→ 要想用 EM 在二模型上得到较好结果  
需满足  $M$  (样本数)  $\gg N$  (样本维数). 这个问题类似于线性方程组求解, 但未知的个数比方程的数目多, 因而不存在完全未知数.

降维方法

→ 假设存在隐含变量  $z \sim N(0, 1)$   $z \in \mathbb{R}^d$  ( $d < n$ )

假设训练样本  $X$  由隐含变量  $z$  生成, 即  $X = \mu + \lambda z + \epsilon$

其中:  $\epsilon \in N(0, \psi)$ . 等价于,  $z$  已知时  $X$  的概率分布为:  $X|z \sim N(\mu + \lambda z, \psi)$  这就是因子分析模型的含义. 其中:  $\mu \in \mathbb{R}^n$ ,  $\lambda \in \mathbb{R}^{n \times d}$ ,  $\psi \in \mathbb{R}^{n \times n}$  是对角矩阵.

主成分分析 (PCA): 将冗余属性去掉: 为方便计算需预处理

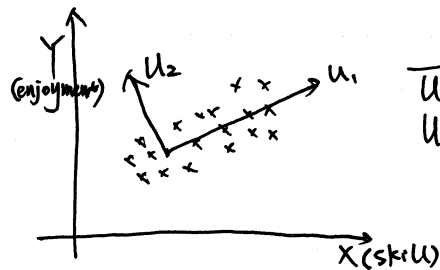
1.  $\mu = \frac{1}{m} \sum_{i=1}^m x^{(i)}$
2.  $x^{(i)} - \mu$  替换  $x^{(i)}$
3.  $\sigma_j^2 = \frac{1}{m} \sum_{i=1}^m (x_j^{(i)})^2$
4.  $x_j^{(i)} / \sigma_j$  替换  $x_j^{(i)}$
- 1-2 使均值为 0
- 3-4 使每个维度的方差为 1, 归一化每个维度的尺度

寻找主方向, 使点在主方向上的投影点方差最大化,

然后忽略掉主方向上的影响 (降维), 继续寻找主方向  $u$

最终得到  $n$  个主方向  $u$  (相互正交). 选取方差最大的  $k$  个  $u$ , 则新的样本为 (降维后)

$$y^{(i)} = \begin{bmatrix} u_1^T x^{(i)} \\ u_2^T x^{(i)} \\ \vdots \\ u_k^T x^{(i)} \end{bmatrix}$$



$x_1, y$  有强相关性.

$u_1$  展拓相关性  
 $u_2$  则是主方向之外的噪声

PCA 的一个应用: 隐含语义索引 LSI

PCA 的一种实现: 奇异值分解 SVD

	密度估计法 (概率方法)	非概率方法
降维到子空间	因子分析	PCA
假设数据位于区块中	混合高斯模型	k-Means

有监督学习：告诉算法每个样品的正确答案，学习后的算法，对新输入也能输出正确答案。

1-8

1. 回归问题：需要预测的变量是连续的

2. 分类问题：需要处理的问题是离散的

用最小二乘法，做契合度评估

求解方法：梯度下降  
正规方程组

$$\theta = (X^T X)^{-1} X^T Y$$

无监督学习：

9-11

根据样本集大小分为：

A. 过拟合 (Overfitting)，特征集过大，

使用特征选择算法，类自动化算法，选取用到的特征

B. 欠拟合 (underfitting)，特征集过小。

使用局部加权回归，来缓解对于选取特征的需求

参数学习算法

参数学习算法：有固定数目参数的数据拟合算法

参数数目会随着训练集增长而变化的算法，即便通过学习之后

学习(型)算法的理论基础

9-11

强化学习：

1. 回报函数 (正反馈/负反馈)

回归算法：logistic  $\Rightarrow$  (分类算法，回归拟合)

$$L(\theta) = PC(Y|X; \theta) \quad \text{似然性}$$

$$l(\theta) = \log L(\theta) \quad \text{对数似然性}$$

使用梯度下降方法，使求似然性最大

$$\theta_j := \theta_j + \alpha \nabla_{\theta} l(\theta)$$

因为求最大，所以梯度上升

偏导数后得到

$$\theta_j := \theta_j + \alpha (y^{(i)} - h(x^{(i)})) x_j^{(i)}$$

$$\nabla_{\theta} l(\theta) = \frac{\partial}{\partial \theta} l(\theta)$$

以上类似于梯度上升算法，形式上和线性回归是相同的，只是符号相反。但

$h(x^{(i)})$  是 logistic 函数，和线性回归的实质是不一样的

采用  $g(z) = \frac{1}{1 + e^{-z}}$   
即 sigmoid 函数

每一次预测，都要重新拟合一条曲线

权重  $w^{(i)} = \exp(-\frac{(x^{(i)} - x)^2}{2\tau^2})$ ，即越靠近  $x^{(i)}$

越接近要预测的输入  $x$ ， $w^{(i)}$  的值越接近 1。

3. 越接近要预测的输入  $x$ ， $w^{(i)}$  的值越接近 1。

4. 若训练集很大，则用于预测的训练集就很大，工作量就呵呵了，此时要用 KD-tree 了

以上可以用概率来解释：

1. 线性回归中，为什么选择最小二乘法作为计算参数的标准，使得假设预测值与真实值之间的面积平方最小？

2. 极大似然估计：选择  $\theta$  使似然性  $L(\theta)$  最大 (数据出现的概率最大)，则要使  $\sum (y^{(i)} - \theta^T x^{(i)})^2$  最小，即之前的  $J(\theta)$ 。所以可知：之前的最小二乘法计算参数，实际上假设了误差项满足高斯分布，且独立分布的情况，使似然最大化来计算参数。

牛顿方法，可以应用于 logistic 回归模型，加速拟合过程。

对于:  $P(y|x; \theta)$

如果,  $y$  为实数, 满足高斯分布, 得到基于最小二乘法的线性回归

如果  $y$  为  $\{0, 1\}$ , 满足伯努利分布, 得到 logistic 回归.

选定了一个指数分布族后, 怎么来推导出一个 GLM (广义线性模型) 呢?



生成学习算法

高斯判别分析 GDA

朴素贝叶斯

Laplace平滑

sigmoid函数

→ 阈值函数

计算  $\theta^T x$ , 输出  $1 \Leftrightarrow \theta^T x > 0$ ,  $0 \Leftrightarrow \theta^T x < 0$ 对于如果  $\theta^T x \gg 0$ , 相当于确定的预测  $y=1$ ,如果  $\theta^T x \ll 0$ , 相当于确定的预测  $y=0$ 对于所有  $i$ , 如果  $y_i=1$ ,  $\theta^T(x^{(i)}) \gg 0$ , 如果  $y_i=0$ ,  $\theta^T(x^{(i)}) \ll 0$ , 则我们  
认为分类器是良好的线性分类方法: Logistic回归 [通过牛顿法或梯度下降得到一条直线, 分割样本]↑  
指数分布族 [联合分布是 log'stic 函数]

非线性分类方法: 神经网络

支持向量机  $\Rightarrow$  线性分类器, 可以拓展为非线性

(函数间隔、几何间隔)?

超平面  $(w, b)$  和某特定训练样本  $(x^{(i)}, y^{(i)})$  的函数间隔定义为:  $\hat{y}^{(i)} = y^{(i)}(w^T x + b)$  $w^T x + b = 0$  → 表示分类的超平面和正确性正确时:  $|w^T x + b|$  表示点到面的距离

→ 法向量

非凸性约束

$$\max \frac{1}{\|w\|} \Leftrightarrow \min \frac{1}{2} \|w\|^2$$

拉格朗日算子  
拉格朗日函数  
对偶问题

$$ax + b + y = 0$$

$$wx + b$$

$$\boxed{\vec{b} + A} \boxed{b} = 0$$