

第四章 概率方法

张皓

<https://haomood.github.io/homepage/>
zhangh0214@gmail.com

摘要

本文介绍机器学习算法中的概率方法。概率方法会对数据的分布进行假设,对概率密度函数进行估计,并使用这个概率密度函数进行决策。本文介绍四种最常用的概率方法:线性回归(用于回归任务)、对数几率回归(用于二分类任务)、*Softmax* 回归(用于多分类任务)和朴素贝叶斯分类器(用于多分类任务)*。前三种方法属于判别式模型,而朴素贝叶斯分类器属于生成式模型。

本系列文章有以下特点:(a). 为了减轻读者的负担并能使尽可能多的读者从中收益,本文试图尽可能少地使用数学知识,只要求读者有基本的微积分、线性代数和概率论基础,并在第一节对关键的数学知识进行回顾和介绍。(b). 本文不省略任何推导步骤,适时补充背景知识,力图使本节内容是自足的,使机器学习的初学者也能理解本文内容。(c). 机器学习近年来发展极其迅速,已成为一个非常广袤的领域。本文无法涵盖机器学习领域的方方面面,仅就一些关键的机器学习流派的方法进行介绍。(d). 为了帮助读者巩固本文内容,或引导读者扩展相关知识,文中穿插了许多问题,并在最后一节进行问题的“快问快答”。

*严格来说,前三者兼有多种解释,既可以看做是概率方法,又可以看做是非概率方法。

1 准备知识

本节给出概率方法的基本流程,后续要介绍的不同概率方法都遵循这一基本流程。

1.1 概率方法的建模流程

(1). 对 $p(y | x, \theta)$ 进行概率假设。我们假定 $p(y | x, \theta)$ 具有某种确定的概率分布形式,其形式被参数向量 θ 唯一地确定。

(2). 对参数 θ 进行最大后验估计。基于训练样例对概率分布的参数 θ 进行最大后验估计 (maximum a posteriori, MAP), 得到需要优化的损失函数。

最大后验估计是指

$$\theta^* := \arg \max_{\theta} p(\theta | D) = \arg \max_{\theta} p(D | \theta) p(\theta), \quad (1)$$

其在最大化时考虑如下两项:

- 参数的先验分布 $p(\theta)$ 。最大后验估计认为参数 θ 未知并且是一个随机变量,其本身服从一个先验分布 $p(\theta)$ 。这个先验分布蕴含了我们关于参数的领域知识。
- 基于观测数据得到的似然 (likelihood) $p(D | \theta)$ 。最大化似然是在 θ 的所有可能的取值中,找到一个能使样本属于其真实标记的概率最大的值。

最大后验估计是在考虑先验分布 $p(\theta)$ 时最大化基于观测数据得到的似然 (likelihood) $p(D | \theta)$ 。

参数估计的两个不同学派的基本观点是什么? 这实际上是参数估计 (parameter estimation) 过程,统计学中的频率主义学派 (frequentist) 和贝叶斯学派 (Bayesian) 提供了不同的解决方案 [3, 9]。频率主义学派认为参数虽然未知,但确实客观存在的固定值,因此通常使用极大似然估计来确定参数值。贝叶斯学派则

认为参数是未观察到的随机变量, 其本身也可有分布, 因此, 可假定参数服从一个先验分布, 然后基于观察到的数据来计算参数的后验分布.

定理 1. 最大后验估计的结果是优化如下形式的损失函数

$$\mathcal{L}(\theta) := \frac{1}{m} \sum_{i=1}^m \ell(\theta) + \lambda \Omega(\theta). \quad (2)$$

Proof. 利用样例的独立同分布假设,

$$\begin{aligned} \theta^* &:= \arg \max_{\theta} p(\theta | D) \\ &= \arg \max_{\theta} p(D | \theta) p(\theta) \\ &= \arg \max_{\theta} \log p(D | \theta) + \log p(\theta) \\ &= \arg \max_{\theta} \log \prod_{i=1}^m p(\mathbf{x}_i, y_i | \theta) + \log p(\theta) \\ &= \arg \min_{\theta} - \sum_{i=1}^m \log p(y_i | \mathbf{x}_i, \theta) - \log p(\theta). \quad (3) \end{aligned}$$

之后, 从 $-\sum_{i=1}^m \log p(y_i | \mathbf{x}_i, \theta)$ 得到 $\frac{1}{m} \sum_{i=1}^m \ell(\theta)$, 从 $-\log p(\theta)$ 得到 $\lambda \Omega(\theta)$. \square

经验风险和结构风险的含义? $\mathcal{L}(\theta)$ 的第一项称为经验风险 (empirical risk), 用于描述模型与训练数据的契合程度. 第二项称为结构风险 (structural risk) 或正则化项 (regularization term), 源于模型的先验概率, 表述了我们希望获得何种性质的模型 (例如希望获得复杂度较小的模型). λ 称为正则化常数, 对两者进行折中.

结构风险的作用? (1). 为引入领域知识和用户意图提供了途径. (2). 有助于削减假设空间, 从而降低了最小化训练误差的过拟合风险. 这也可理解为一种“罚函数法”, 即对不希望得到的结果施以惩罚, 从而使得优化过程趋向于希望目标. ℓ_p 范数是常用的正则化项.

定理 2. 若参数 θ 的先验分布服从高斯分布 $\theta \sim \mathcal{N}(\mathbf{0}, \sigma^2 \mathbf{I})$, 那么对应的结构风险是 $\Omega(\theta) := \frac{1}{2} \|\theta\|^2$.

Proof. 由于从 $-\log p(\theta)$ 得到 $\lambda \Omega(\theta)$,

$$\begin{aligned} \theta^* &= \arg \min_{\theta} -\log p(\theta) \\ &= \arg \min_{\theta} -\log \frac{1}{\sqrt{(2\pi)^{d+1} \det \sigma^2 \mathbf{I}}} \exp\left(-\frac{1}{2} \theta^\top \sigma^2 \mathbf{I} \theta\right) \\ &= \arg \min_{\theta} \frac{1}{2} \theta^\top \sigma^2 \mathbf{I} \theta \\ &= \arg \min_{\theta} \frac{1}{2} \|\theta\|^2, \quad (4) \end{aligned}$$

其中先验分布 $\theta \sim \mathcal{N}(\mathbf{0}, \sigma^2 \mathbf{I})$ 的参数 σ^2 转化为正则化常数 λ . \square

为什么最常假设参数的先验分布是高斯分布 (或最常使用 ℓ_2 正则化)? 这是因为高斯分布 $\mathcal{N}(\mu, \Sigma)$ 是所有均值和熵存在且协方差矩阵是 Σ 的分布中熵最大的分布. 最大熵分布是在特定约束下具有最大不确定性的分布. 在没有更多信息的情况下, 那些不确定的部分都是“等可能的”. 在设计先验分布 $p(\theta)$ 时, 除了我们对参数的认知 (例如均值和值域) 外, 我们不想引入任何其余的偏见 (bias). 因此最大熵先验 (对应 ℓ_2 正则化) 常被使用. 除高斯先验外, 还可以使用不提供信息的先验 (uninformative prior), 其在一定范围内均匀分布, 对应的损失函数中没有结构风险这一项.

(3). 对损失函数 $\mathcal{L}(\theta)$ 进行梯度下降优化.

$$\theta \leftarrow \theta - \eta \frac{\partial \mathcal{L}}{\partial \theta}. \quad (5)$$

梯度下降的细节留在下一节介绍.

概率方法的优缺点各是什么? 优点: 这种参数化的概率方法使参数估计变得相对简单. 缺点: 参数估计结果的准确性严重依赖于所假设的概率分布形式是否符合潜在的真实数据分布. 在现实应用中, 欲做出能较好地接近潜在真实分布的假设, 往往需在一定程度利用关于应用任务本身的经验知识, 否则仅凭“猜测”来假设概率分布形式, 很可能产生误导性的结果. 我们不一定非要概率式地解释这个世界, 在不考虑概率的情况下, 直接找到分类边界, 也被称为判别函数 (discriminant function) 有时甚至能比判别式模型产生更好的结果.

1.2 梯度下降

我们的目标是求解下列无约束的优化问题

$$\theta^* := \arg \min_{\theta} \mathcal{L}(\theta), \quad (6)$$

其中 $\mathcal{L}(\theta)$ 是连续可微函数. 梯度下降是一种一阶 (first-order) 优化方法, 是求解无约束优化问题最简单、最经典的求解方法之一.

梯度下降的基本思路? 梯度下降贪心地迭代式地最小化 $\mathcal{L}(\theta)$. 梯度下降希望找到一个方向 (单位向量) \mathbf{v} 使得 \mathcal{L} 在这个方向下降最快, 并在这个方向前进 α 的距离

$$\theta \leftarrow \theta + \alpha \mathbf{v}. \quad (7)$$

定理 3. 梯度下降的更新规则是公式 5. 重复这个过程, 可收敛到局部极小点.

Proof. 我们需要找到下降最快的方向 \mathbf{v} 和前进的距离 α .

(1). 下降最快的方向 \mathbf{v} . 利用泰勒展开

$$f(\mathbf{x}) \approx f(\mathbf{x}_0) + (\mathbf{x} - \mathbf{x}_0)^\top \frac{\partial f}{\partial \mathbf{x}} \Big|_{\mathbf{x}_0} + \frac{1}{2} (\mathbf{x} - \mathbf{x}_0)^\top \frac{\partial^2 f}{\partial \mathbf{x}^2} \Big|_{\mathbf{x}_0} (\mathbf{x} - \mathbf{x}_0) \quad (8)$$

的一阶近似,

$$\begin{aligned} \mathbf{v}^* &:= \arg \min_{\|\mathbf{v}\|=1} \mathcal{L}(\boldsymbol{\theta} + \alpha \mathbf{v}) \\ &\approx \arg \min_{\|\mathbf{v}\|=1} \mathcal{L}(\boldsymbol{\theta}) + \alpha \mathbf{v}^\top \frac{\partial \mathcal{L}}{\partial \boldsymbol{\theta}} \\ &= \arg \min_{\|\mathbf{v}\|=1} \alpha \|\mathbf{v}\| \left\| \frac{\partial \mathcal{L}}{\partial \boldsymbol{\theta}} \right\| \cos \left\langle \mathbf{v}, \frac{\partial \mathcal{L}}{\partial \boldsymbol{\theta}} \right\rangle \\ &= \arg \min_{\|\mathbf{v}\|=1} \cos \left\langle \mathbf{v}, \frac{\partial \mathcal{L}}{\partial \boldsymbol{\theta}} \right\rangle \\ &= - \frac{\frac{\partial \mathcal{L}}{\partial \boldsymbol{\theta}}}{\left\| \frac{\partial \mathcal{L}}{\partial \boldsymbol{\theta}} \right\|}, \end{aligned} \quad (9)$$

即下降最快的方向是损失函数的负梯度方向.

(2). 前进的距离 α . 我们希望在开始的时候前进距离大一些以使得收敛比较快, 而在接近最小值时前进距离小一些以不错过最小值点. 因此, 我们设前进距离为损失函数梯度的一个倍数

$$\alpha^* := \eta \left\| \frac{\partial \mathcal{L}}{\partial \boldsymbol{\theta}} \right\|, \quad (10)$$

其中 η 被称为学习率 (learning rate).

向公式 7 代入最优的 α^* 和 \mathbf{v}^* 后即得. \square

引理 4. 若 $\mathcal{L}(\boldsymbol{\theta})$ 满足 L -Lipschitz 条件, 则设置步长 $\eta := \frac{1}{2L}$ 即可确保收敛到局部极小点.

定义 1 (凸函数). 对区间 $[a, b]$ 上定义的函数 $f: \mathbb{R}^n \rightarrow \mathbb{R}$, 若它对区间中任意两点 $\mathbf{x}_1, \mathbf{x}_2$ 均有

$$f\left(\frac{\mathbf{x}_1 + \mathbf{x}_2}{2}\right) \leq \frac{f(\mathbf{x}_1) + f(\mathbf{x}_2)}{2}, \quad (11)$$

则称 f 为区间 $[a, b]$ 上的凸函数 (convex function). 当 $<$ 成立时, 称为严格凸函数 (strict convex function). U 形曲线的函数如 $f(x) := x^2$ 通常是凸函数.

引理 5 (凸函数判定准则). 对 $f: \mathbb{R} \rightarrow \mathbb{R}$, 若 $\frac{d^2 f}{dx^2} \geq 0 (> 0)$, 则 f 是凸函数 (严格凸函数). 对 $f: \mathbb{R}^d \rightarrow \mathbb{R}$, 若 $\frac{\partial^2 f}{\partial \mathbf{x}^2} \succeq 0 (> 0)$, 即 $\frac{\partial^2 f}{\partial \mathbf{x}^2}$ 是一个半正定 (正定) 矩阵, 则 f 是凸函数 (严格凸函数).

引理 6. 对凸函数 $f: \mathbb{R}^d \rightarrow \mathbb{R}$, 局部极小点对应着函数的全局最小点, 此时梯度下降可确保收敛到全局最优解. 全局最小点满足 $\frac{\partial f}{\partial \mathbf{x}} = \mathbf{0}$ 的 \mathbf{x} .

2 线性回归

2.1 建模流程

线性回归 (linear regression) 回归问题 $y \in \mathbb{R}$. 其建模方法包括如下三步 (参见第 1.1 节).

(1). 对 $p(y | \mathbf{x}, \boldsymbol{\theta})$ 进行概率假设.

我们假设

$$y := \mathbf{w}^\top \mathbf{x} + b + \varepsilon, \quad \varepsilon \sim \mathcal{N}(0, \sigma^2). \quad (12)$$

ε 被称为误差项, 捕获了 (a). 特征向量 \mathbf{x} 中没有包含的因素. (b). 随机噪声. 对不同的样本 \mathbf{x}_i , ε_i 是独立同分布地从 $\mathcal{N}(0, \sigma^2)$ 中进行采样得到的.

线性回归的假设函数是

$$h(\mathbf{x}) := \mathbf{w}^\top \mathbf{x} + b. \quad (13)$$

为了书写方便, 我们记

$$\boldsymbol{\theta} := \begin{bmatrix} \mathbf{w} \\ b \end{bmatrix}, \quad \tilde{\mathbf{x}} := \begin{bmatrix} \mathbf{x} \\ 1 \end{bmatrix}. \quad (14)$$

那么公式 12 等价于

$$y = \boldsymbol{\theta}^\top \tilde{\mathbf{x}} + \varepsilon, \quad \varepsilon \sim \mathcal{N}(0, \sigma^2). \quad (15)$$

在本文其余部分我们将沿用这一简化记号. 因此,

$$p(y | \mathbf{x}, \boldsymbol{\theta}) = \mathcal{N}(\boldsymbol{\theta}^\top \tilde{\mathbf{x}}, \sigma^2). \quad (16)$$

(2). 对参数 $\boldsymbol{\theta}$ 进行最大后验估计.

定理 7. 假设参数 $\boldsymbol{\theta}$ 服从高斯先验, 对参数 $\boldsymbol{\theta}$ 进行最大后验估计等价于优化如下损失函数

$$\mathcal{L}(\boldsymbol{\theta}) := \frac{1}{2m} \sum_{i=1}^m (\boldsymbol{\theta}^\top \tilde{\mathbf{x}}_i - y_i)^2 + \lambda \|\boldsymbol{\theta}\|^2. \quad (17)$$

其中

$$\ell(\boldsymbol{\theta}) := \frac{1}{2} (\boldsymbol{\theta}^\top \tilde{\mathbf{x}} - y)^2 \quad (18)$$

被称为平方损失 (square loss). 在线性回归中, 平方损失就是试图找到一个超平面 $\mathbf{w}^\top \mathbf{x} + b = 0$, 使所有样本到该超平面的欧式距离 (Euclidean distance) 之和最小.

Table 1: 常用线性回归模型对比.

学习算法	噪声分布	参数先验	损失函数	特点
简单线性回归	高斯	无信息先验	$\frac{1}{2m} \sum_{i=1}^m (\boldsymbol{\theta}^\top \mathbf{x}_i - y_i)^2$	当属性数 d 大, 样本数 m 小时, 容易过拟合
岭回归 [12]	高斯	高斯	$\frac{1}{2m} \sum_{i=1}^m (\boldsymbol{\theta}^\top \mathbf{x}_i - y_i)^2 + \frac{\lambda}{2} \ \boldsymbol{\theta}\ ^2$	相比简单线性回归, 可以显著降低过拟合风险
LASSO [10]	高斯	拉普拉斯	$\frac{1}{2m} \sum_{i=1}^m (\boldsymbol{\theta}^\top \mathbf{x}_i - y_i)^2 + \lambda \ \boldsymbol{\theta}\ _1$	降低过拟合风险, 并更容易得到稀疏解
稳健线性回归	拉普拉斯	高斯	$\frac{1}{m} \sum_{i=1}^m \boldsymbol{\theta}^\top \mathbf{x}_i - y_i + \frac{\lambda}{2} \ \boldsymbol{\theta}\ ^2$	适用于数据存在异常点

Proof.

$$\begin{aligned}
\boldsymbol{\theta}^* &= \arg \min_{\boldsymbol{\theta}} - \sum_{i=1}^m \log p(y_i | \mathbf{x}_i, \boldsymbol{\theta}) + \lambda \|\boldsymbol{\theta}\|^2 \\
&= \arg \min_{\boldsymbol{\theta}} - \sum_{i=1}^m \log \frac{1}{\sqrt{2\pi\sigma^2}} \exp\left(-\frac{(y_i - \boldsymbol{\theta}^\top \tilde{\mathbf{x}}_i)^2}{2\sigma^2}\right) + \lambda \|\boldsymbol{\theta}\|^2 \\
&= \arg \min_{\boldsymbol{\theta}} \sum_{i=1}^m (y_i - \boldsymbol{\theta}^\top \tilde{\mathbf{x}}_i)^2 + \lambda \|\boldsymbol{\theta}\|^2 \\
&= \arg \min_{\boldsymbol{\theta}} \frac{1}{2m} \sum_{i=1}^m (\boldsymbol{\theta}^\top \tilde{\mathbf{x}}_i - y_i)^2 + \lambda \|\boldsymbol{\theta}\|^2.
\end{aligned}$$

其中, 最后一行只是为了数学计算上方便, 下文推导对数几率回归和 Softmax 回归时的最后一步亦然. \square

(3). 对损失函数 $\mathcal{L}(\boldsymbol{\theta})$ 进行梯度下降优化.

可以容易地得到损失函数对参数的偏导数

$$\frac{\partial \mathcal{L}}{\partial \boldsymbol{\theta}} = \frac{1}{m} \sum_{i=1}^m (\boldsymbol{\theta}^\top \tilde{\mathbf{x}}_i - y_i) \tilde{\mathbf{x}}_i + \lambda \boldsymbol{\theta}. \quad (20)$$

2.2 线性回归的闭式解

线性回归对应的平方损失的函数形式比较简单, 可以通过求 $\frac{\partial \mathcal{L}}{\partial \boldsymbol{\theta}} = \mathbf{0}$ 直接得到最优解.

定理 8. 线性回归的闭式解为

$$\boldsymbol{\theta}^* = (\mathbf{X}^\top \mathbf{X} + m\lambda \mathbf{I})^{-1} \mathbf{X}^\top \mathbf{y}. \quad (21)$$

Proof. $\mathcal{L}(\boldsymbol{\theta})$ 可等价地写作

$$\begin{aligned}
\mathcal{L}(\boldsymbol{\theta}) &= \frac{1}{2m} \sum_{i=1}^m (\boldsymbol{\theta}^\top \tilde{\mathbf{x}}_i - y_i)^2 + \frac{\lambda}{2} \|\boldsymbol{\theta}\|^2 \\
&= \frac{1}{2m} \sum_{i=1}^m (\tilde{\mathbf{x}}_i^\top \boldsymbol{\theta} - y_i)^2 + \frac{\lambda}{2} \boldsymbol{\theta}^\top \boldsymbol{\theta} \\
&= \frac{1}{2m} \left\| \begin{bmatrix} \tilde{\mathbf{x}}_1^\top \boldsymbol{\theta} - y_1 \\ \tilde{\mathbf{x}}_2^\top \boldsymbol{\theta} - y_2 \\ \vdots \\ \tilde{\mathbf{x}}_m^\top \boldsymbol{\theta} - y_m \end{bmatrix} \right\|^2 + \frac{\lambda}{2} \boldsymbol{\theta}^\top \boldsymbol{\theta}. \quad (22)
\end{aligned}$$

令

$$\mathbf{X} := \begin{bmatrix} \tilde{\mathbf{x}}_1^\top \\ \tilde{\mathbf{x}}_2^\top \\ \vdots \\ \tilde{\mathbf{x}}_m^\top \end{bmatrix} \in \mathbb{R}^{m \times (d+1)}, \quad \mathbf{y} := \begin{bmatrix} y_1 \\ y_2 \\ \vdots \\ y_m \end{bmatrix} \in \mathbb{R}^m, \quad (23)$$

那么

$$\begin{aligned}
(19) \quad \mathcal{L}(\boldsymbol{\theta}) &= \frac{1}{2m} \|\mathbf{X}\boldsymbol{\theta} - \mathbf{y}\|^2 + \frac{\lambda}{2} \boldsymbol{\theta}^\top \boldsymbol{\theta} \\
&= \frac{1}{2m} (\boldsymbol{\theta}^\top \mathbf{X}^\top \mathbf{X} \boldsymbol{\theta} - 2\boldsymbol{\theta}^\top \mathbf{X}^\top \mathbf{y} + \mathbf{y}^\top \mathbf{y}) + \frac{\lambda}{2} \boldsymbol{\theta}^\top \boldsymbol{\theta} \quad (24)
\end{aligned}$$

求解

$$\frac{\partial \mathcal{L}}{\partial \boldsymbol{\theta}} = \frac{1}{m} \mathbf{X}^\top \mathbf{X} \boldsymbol{\theta} - \frac{1}{m} \mathbf{X}^\top \mathbf{y} + \lambda \boldsymbol{\theta} = \mathbf{0} \quad (25)$$

即得. \square

$\mathbf{X}^\top \mathbf{X}$ 不可逆的情况及解决方案? (1). 属性数 $d+1$ 多于样例数 m . (2). 属性之间线性相关. 通过正则化项 $m\lambda \mathbf{I}$, 即使 $\mathbf{X}^\top \mathbf{X}$ 不可逆, $\mathbf{X}^\top \mathbf{X} + m\lambda \mathbf{I}$ 仍是可逆的.

2.3 其他正则化回归模型

事实上, 上文介绍的线性回归模型是岭回归 (ridge regression). 根据正则化项的不同, 有三种常用的线性回归模型, 见表 1.

基于 ℓ_0 、 ℓ_1 和 ℓ_2 范数正则化的效果? ℓ_2 范数倾向于 \mathbf{w} 的分量取值尽量均衡, 即非零分量个数尽量稠密. 而 ℓ_0 “范数”和 ℓ_1 范数则倾向于 \mathbf{w} 的分量尽量稀疏, 即非零分量个数尽量少, 优化结果得到了仅采用一部分属性的模型. 也就是说, 基于 ℓ_0 “范数”和 ℓ_1 范数正则化的学习方法是一种嵌入式 (embedding) 特征选择方法, 其特征选择过程和学习器训练过程融为一体, 两者在同一个优化过程中完成. 事实上, 对 \mathbf{w} 施加稀疏约束最自

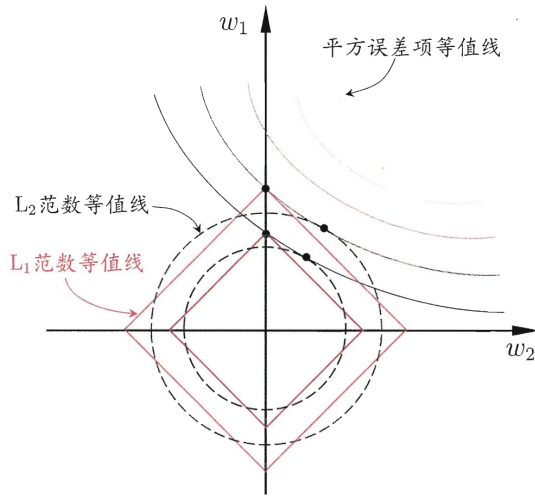


Figure 1: ℓ_1 正则化 (红色) 比 ℓ_2 正则化 (黑色) 更易于获得稀疏解. 本图源于 [17].

然的是使用 ℓ_0 “范数”. 但 ℓ_0 “范数”不连续, 难以优化求解. 因此常采用 ℓ_1 范数来近似.

为什么 ℓ_1 正则化比 ℓ_2 正则化更易于获得稀疏解? 假设 $\mathbf{x} \in \mathbb{R}^2$, 则 $\mathbf{w} \in \mathbb{R}^2$. 我们绘制出平方损失项、 ℓ_1 范数和 ℓ_2 范数的等值线 (取值相同的点的连线), 如图 1 所示. LASSO 的解要在平方损失项和正则化项之间折中, 即出现在图中平方误差项等值线和正则化项等值线的相交处. 从图中可以看出, 采用 ℓ_1 正则化时交点常出现在坐标轴上 ($w_2 = 0$), 而采用 ℓ_2 正则化时交点常出现在某个象限中 (w_1, w_2 均不为 0).

考虑一般的带有 ℓ_1 正则化的优化目标

$$\boldsymbol{\theta}^* = \arg \min_{\boldsymbol{\theta}} \ell(\boldsymbol{\theta}) + \lambda \|\boldsymbol{\theta}\|_1, \quad (26)$$

若 $\ell(\boldsymbol{\theta})$ 满足 L -Lipschitz 条件, 即

$$\exists L > 0, \forall \mathbf{x}_1, \mathbf{x}_2, \left\| \frac{\partial \ell}{\partial \mathbf{x}} \Big|_{\mathbf{x}_1} - \frac{\partial \ell}{\partial \mathbf{x}} \Big|_{\mathbf{x}_2} \right\|^2 \leq L \|\mathbf{x}_1 - \mathbf{x}_2\|^2. \quad (27)$$

优化通常使用近端梯度下降 (proximal gradient descent, PGD) [1]. PGD 也是一种贪心地迭代式地最小化策略, 能快速求解基于 ℓ_1 范数最小化的方法.

定理 9. 假设当前参数是 $\boldsymbol{\theta}_t$, PGD 的更新准则是

$$\boldsymbol{\theta}_{t+1,j} := \begin{cases} z_j + \frac{\lambda}{L} & \text{若 } z_j < -\frac{\lambda}{L}; \\ 0 & \text{若 } -\frac{\lambda}{L} \leq z_j \leq \frac{\lambda}{L}; \\ z_j - \frac{\lambda}{L} & \text{若 } z_j > \frac{\lambda}{L}, \end{cases} \quad (28)$$

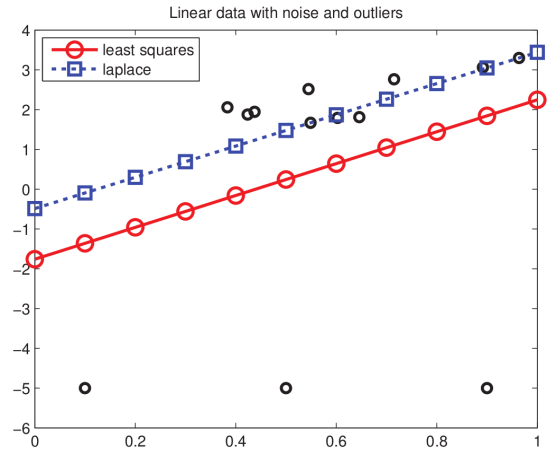


Figure 2: 存在异常点 (图下方的三个点) 时普通线性回归 (红色) 和稳健线性回归 (蓝色). 本图源于 [7].

其中

$$\mathbf{z} := \boldsymbol{\theta}_t - \frac{1}{L} \frac{\partial \ell}{\partial \boldsymbol{\theta}} \Big|_{\boldsymbol{\theta}_t}. \quad (29)$$

Proof. 在 $\boldsymbol{\theta}_t$ 附近将 $\ell(\boldsymbol{\theta})$ 进行二阶泰勒展开近似

$$\begin{aligned} \boldsymbol{\theta}_{t+1} &= \arg \min_{\boldsymbol{\theta}} \ell(\boldsymbol{\theta}) + \lambda \|\boldsymbol{\theta}\|_1 \\ &= \arg \min_{\boldsymbol{\theta}} \ell(\boldsymbol{\theta}_t) + (\boldsymbol{\theta} - \boldsymbol{\theta}_t)^\top \frac{\partial \ell}{\partial \boldsymbol{\theta}} \Big|_{\boldsymbol{\theta}_t} + \frac{L}{2} \|\boldsymbol{\theta} - \boldsymbol{\theta}_t\|^2 \\ &\quad + \lambda \|\boldsymbol{\theta}\|_1 \\ &= \arg \min_{\boldsymbol{\theta}} \frac{L}{2} \left\| \boldsymbol{\theta} - \left(\boldsymbol{\theta}_t - \frac{1}{L} \frac{\partial \ell}{\partial \boldsymbol{\theta}} \Big|_{\boldsymbol{\theta}_t} \right) \right\|^2 + \lambda \|\boldsymbol{\theta}\|_1 \\ &= \arg \min_{\boldsymbol{\theta}} \sum_{j=1}^d \left(\frac{1}{2} (\theta_d - z_d)^2 + \frac{\lambda}{L} |\theta_d| \right). \end{aligned} \quad (30)$$

由于 $\boldsymbol{\theta}$ 各维互不影响 (不存在交叉项), 因此可以独立求解各维. \square

在 LASSO 的基础上进一步发展出考虑特征分组的 Group LASSO [14]、考虑特征序结构的 Fused LASSO [11] 等变体. 由于凸性不严格, LASSO 类方法可能产生多个解, 该问题通过弹性网 (elastic net) 得以解决 [16].

2.4 存在异常点数据的线性回归

一旦数据中存在异常点 (outlier), 由于平方损失计算的是样本点到超平面距离的平方, 远离超平面的点对回归结果产生更大的影响, 如图 2 所示. 平方损失对应于假设噪声服从高斯分布 $\varepsilon \sim \mathcal{N}(0, \sigma^2)$, 一种应对

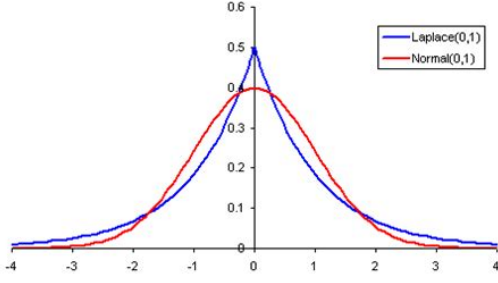


Figure 3: 高斯分布 $\mathcal{N}(0, 1)$ (红色) 和拉普拉斯分布 $\text{Lap}(0, 1)$ (蓝色). 本图源于 <https://www.epixanalytics.com/modelassist/AtRisk/images/15/image632.gif>.

异常点的方法是取代高斯分布为其他更加重尾 (heavy tail) 的分布, 使其对异常点的容忍能力更强, 例如使用拉普拉斯分布 $\varepsilon \sim \text{Lap}(0, b)$, 如图 3 所示.

定义 2 (拉普拉斯分布 (Laplace distribution) $\text{Lap}(\mu, b)$). 又称为双边指数分布 (double sided exponential distribution), 具有如下的概率密度函数

$$p(x) := \frac{1}{2b} \exp\left(-\frac{|x - \mu|}{b}\right). \quad (31)$$

该分布均值为 μ , 方差为 $2b^2$.

定理 10. 假设参数服从高斯先验,

$$p(y | \mathbf{x}, \boldsymbol{\theta}) = \text{Lap}(\boldsymbol{\theta}^\top \tilde{\mathbf{x}}, b), \quad (32)$$

对参数 $\boldsymbol{\theta}$ 进行最大后验估计等价于优化如下损失函数

$$\mathcal{L}(\boldsymbol{\theta}) := \frac{1}{m} \sum_{i=1}^m |\boldsymbol{\theta}^\top \tilde{\mathbf{x}}_i - y_i| + \lambda \|\boldsymbol{\theta}\|^2. \quad (33)$$

Proof.

$$\begin{aligned} \boldsymbol{\theta}^* &= \arg \min_{\boldsymbol{\theta}} - \sum_{i=1}^m \log p(y_i | \mathbf{x}_i, \boldsymbol{\theta}) + \lambda \|\boldsymbol{\theta}\|^2 \\ &= \arg \min_{\boldsymbol{\theta}} - \sum_{i=1}^m \log \frac{1}{2b} \exp\left(-\frac{|y_i - \boldsymbol{\theta}^\top \tilde{\mathbf{x}}_i|}{b}\right) + \lambda \|\boldsymbol{\theta}\|^2 \\ &= \arg \min_{\boldsymbol{\theta}} \frac{1}{m} \sum_{i=1}^m |\boldsymbol{\theta}^\top \tilde{\mathbf{x}}_i - y_i| + \lambda \|\boldsymbol{\theta}\|^2. \end{aligned} \quad (34)$$

□

由于绝对值函数不光滑, 不便基于梯度下降对公式 33 进行优化. 通过分离变量技巧, 可将其转化为二次规划 (quadratic programming) 问题, 随后调用现有的

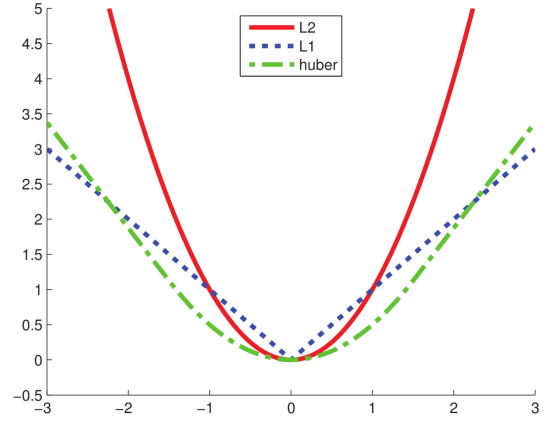


Figure 4: ℓ_2 损失 (红色)、 ℓ_1 损失 (蓝色) 和 Huber 损失 (绿色). 本图源于 [7].

软件包进行求解. 我们在下一章形式化 SVR 时还会再使用这个技巧.

定理 11. 最小化公式 33 等价于如下二次规划问题, 其包含 $d + 1 + 2m$ 个变量, $3m$ 个约束:

$$\begin{aligned} \arg \min_{\boldsymbol{\theta}, \varepsilon^+, \varepsilon^-} & \frac{1}{m} \sum_{i=1}^m (\varepsilon_i^+ - \varepsilon_i^-) + \lambda \|\boldsymbol{\theta}\|^2 \\ \text{s.t.} & \varepsilon_i^+ \geq 0, & i = 1, 2, \dots, m, \\ & \varepsilon_i^- \geq 0, & i = 1, 2, \dots, m, \\ & \varepsilon_i^+ - \varepsilon_i^- + y_i - \boldsymbol{\theta}^\top \tilde{\mathbf{x}}_i = 0 & i = 1, 2, \dots, m. \end{aligned} \quad (35)$$

Proof. 令 $\varepsilon_i^+ - \varepsilon_i^- := \boldsymbol{\theta}^\top \tilde{\mathbf{x}}_i - y_i$, 并约束 $\varepsilon_i^+ \geq 0, \varepsilon_i^- \geq 0$. □

此外, 为了结合高斯分布 (对应平凡损失) 容易优化和拉普拉斯分布 (对应 ℓ_1 损失) 可以应对异常值的优点, Huber 损失 [5] 在误差接近 0 时为平方损失, 在误差比较大时接近 ℓ_1 损失, 如图 4 所示.

$$\ell(\boldsymbol{\theta}) := \begin{cases} \frac{1}{2} (\boldsymbol{\theta}^\top \tilde{\mathbf{x}} - y)^2 & \text{若 } |\boldsymbol{\theta}^\top \tilde{\mathbf{x}} - y| \leq \tau; \\ \tau |\boldsymbol{\theta}^\top \tilde{\mathbf{x}} - y| - \frac{\tau^2}{2} & \text{若 } |\boldsymbol{\theta}^\top \tilde{\mathbf{x}} - y| > \tau. \end{cases} \quad (36)$$

Huber 损失处处可微, 使用基于梯度的方法对 Huber 损失进行优化会比使用拉普拉斯分布更快.

2.5 广义线性模型

线性回归利用属性的线性组合 $\boldsymbol{\theta}^\top \tilde{\mathbf{x}}$ 进行预测. 除了直接利用 $\boldsymbol{\theta}^\top \tilde{\mathbf{x}}$ 逼近 y 外, 还可以使模型的预测值逼

其中

$$\ell(\boldsymbol{\theta}) := -y\boldsymbol{\theta}^\top \tilde{\mathbf{x}} + \log(1 + \exp \boldsymbol{\theta}^\top \tilde{\mathbf{x}}) \quad (44)$$

称为对数几率损失 (logistic loss).

Proof.

$$\begin{aligned} \boldsymbol{\theta}^* &= \arg \min_{\boldsymbol{\theta}} - \sum_{i=1}^m \log p(y_i | \mathbf{x}_i, \boldsymbol{\theta}) + \frac{\lambda}{2} \|\boldsymbol{\theta}\|^2 \\ &= \arg \min_{\boldsymbol{\theta}} - \sum_{i=1}^m \left(y_i \log \text{sigm} \boldsymbol{\theta}^\top \tilde{\mathbf{x}}_i \right. \\ &\quad \left. + (1 - y_i) \log(1 - \text{sigm} \boldsymbol{\theta}^\top \tilde{\mathbf{x}}_i) \right) \\ &\quad + \frac{\lambda}{2} \|\boldsymbol{\theta}\|^2. \end{aligned} \quad (45)$$

注意到

$$\begin{aligned} \log \text{sigm} \boldsymbol{\theta}^\top \tilde{\mathbf{x}}_i &= \log \frac{1}{1 + \exp(-\boldsymbol{\theta}^\top \tilde{\mathbf{x}}_i)} \\ &= \log \frac{\exp \boldsymbol{\theta}^\top \tilde{\mathbf{x}}_i}{1 + \exp \boldsymbol{\theta}^\top \tilde{\mathbf{x}}_i} \\ &= \boldsymbol{\theta}^\top \tilde{\mathbf{x}}_i - \log(1 + \exp \boldsymbol{\theta}^\top \tilde{\mathbf{x}}_i), \end{aligned} \quad (46)$$

$$\begin{aligned} \log(1 - \text{sigm} \boldsymbol{\theta}^\top \tilde{\mathbf{x}}_i) &= \log \left(1 - \frac{1}{1 + \exp(-\boldsymbol{\theta}^\top \tilde{\mathbf{x}}_i)} \right) \\ &= \log \frac{1}{1 + \exp \boldsymbol{\theta}^\top \tilde{\mathbf{x}}_i} \\ &= -\log(1 + \exp \boldsymbol{\theta}^\top \tilde{\mathbf{x}}_i). \end{aligned} \quad (47)$$

因此

$$\begin{aligned} \boldsymbol{\theta}^* &= \arg \min_{\boldsymbol{\theta}} - \sum_{i=1}^m \left(y_i (\boldsymbol{\theta}^\top \tilde{\mathbf{x}}_i - \log(1 + \exp \boldsymbol{\theta}^\top \tilde{\mathbf{x}}_i)) \right. \\ &\quad \left. - (1 - y_i) \log(1 + \exp \boldsymbol{\theta}^\top \tilde{\mathbf{x}}_i) \right) + \frac{\lambda}{2} \|\boldsymbol{\theta}\|^2 \\ &= \arg \min_{\boldsymbol{\theta}} \sum_{i=1}^m \left(-y_i \boldsymbol{\theta}^\top \tilde{\mathbf{x}}_i + \log(1 + \exp \boldsymbol{\theta}^\top \tilde{\mathbf{x}}_i) \right) + \frac{\lambda}{2} \|\boldsymbol{\theta}\|^2 \\ &= \arg \min_{\boldsymbol{\theta}} \frac{1}{m} \sum_{i=1}^m \left(-y_i \boldsymbol{\theta}^\top \tilde{\mathbf{x}}_i + \log(1 + \exp \boldsymbol{\theta}^\top \tilde{\mathbf{x}}_i) \right) \\ &\quad + \frac{\lambda}{2} \|\boldsymbol{\theta}\|^2. \end{aligned} \quad (48)$$

□

(3). 对损失函数 $\mathcal{L}(\boldsymbol{\theta})$ 进行梯度下降优化.

$$\begin{aligned} \frac{\partial \mathcal{L}}{\partial \boldsymbol{\theta}} &= \frac{1}{m} \sum_{i=1}^m \left(-y_i \mathbf{x}_i + \frac{\exp \boldsymbol{\theta}^\top \tilde{\mathbf{x}}_i}{1 + \exp \boldsymbol{\theta}^\top \tilde{\mathbf{x}}_i} \mathbf{x}_i \right) + \lambda \boldsymbol{\theta} \\ &= \frac{1}{m} \sum_{i=1}^m \left(\text{sigm} \boldsymbol{\theta}^\top \tilde{\mathbf{x}}_i - y_i \right) \mathbf{x}_i + \lambda \boldsymbol{\theta}. \end{aligned} \quad (49)$$

3.2 与广义线性模型的关系

对数几率回归的假设函数 $h(\mathbf{x}) = \frac{1}{1 + \exp(-\boldsymbol{\theta}^\top \tilde{\mathbf{x}})}$ 等价于

$$\log \frac{h(\mathbf{x})}{1 - h(\mathbf{x})} = \boldsymbol{\theta}^\top \tilde{\mathbf{x}}, \quad (50)$$

其中

$$\frac{h(\mathbf{x})}{1 - h(\mathbf{x})} = \frac{p(y = 1 | \mathbf{x}, \boldsymbol{\theta})}{p(y = 0 | \mathbf{x}, \boldsymbol{\theta})} \quad (51)$$

被称为几率 (odds), 反映 \mathbf{x} 作为正例的相对可能性. $\log \frac{h(\mathbf{x})}{1 - h(\mathbf{x})}$ 被称为对数几率 (log odds, logit), 公式 50 实际上在用线性回归模型的预测结果逼近真实标记的对数几率, 这是对数几率回归名称的由来.

对数几率回归的优点? (1). 直接对分类的可能性进行建模 (假设 $p(y | \mathbf{x}, \boldsymbol{\theta})$ 服从伯努利分布), 无需事先假设样本 \mathbf{x} 的分布, 这样避免了假设分布不准确所带来的问题. (2). 不仅能预测出类别, 还可以得到近似概率预测, 对许多需要概率辅助决策的任务很有用. (3). 对数几率的目标函数是凸函数, 有很好的数学性质.

引理 13. 对数几率损失函数是凸函数.

Proof. 在 $\frac{\partial \mathcal{L}}{\partial \boldsymbol{\theta}}$ 的基础上, 进一步可求得

$$\frac{\partial^2 \mathcal{L}}{\partial \boldsymbol{\theta}^2} = \frac{1}{m} \sum_{i=1}^m \text{sigm} \boldsymbol{\theta}^\top \tilde{\mathbf{x}}_i (1 - \text{sigm} \boldsymbol{\theta}^\top \tilde{\mathbf{x}}_i) \tilde{\mathbf{x}}_i \tilde{\mathbf{x}}_i^\top + \lambda \quad (52)$$

是一个半正定矩阵. □

3.3 $y \in \{-1, 1\}$ 的对数几率回归

为了概率假设方便, 我们令二分类问题的标记 $y \in \{0, 1\}$. 有时, 我们需要处理 $y \in \{-1, 1\}$ 形式的分类问题. 对数几率损失函数需要进行相应的改动.

(1). 对 $p(y | \mathbf{x}, \boldsymbol{\theta})$ 进行概率假设.

我们假设

$$\Pr(y = 1 | \mathbf{x}, \boldsymbol{\theta}) := \text{sigm} \boldsymbol{\theta}^\top \tilde{\mathbf{x}}, \quad (53)$$

那么

$$\Pr(y = 0 | \mathbf{x}, \boldsymbol{\theta}) = 1 - \text{sigm} \boldsymbol{\theta}^\top \tilde{\mathbf{x}} = \text{sigm}(-\boldsymbol{\theta}^\top \tilde{\mathbf{x}}). \quad (54)$$

两者可以合并写作

$$p(y | \mathbf{x}, \boldsymbol{\theta}) = \text{sigm} y \boldsymbol{\theta}^\top \tilde{\mathbf{x}} \quad (55)$$

(2). 对参数 $\boldsymbol{\theta}$ 进行最大后验估计.

定理 14. 假设参数 θ 服从高斯先验, 对参数 θ 进行最大后验估计等价于优化如下损失函数

$$\mathcal{L}(\theta) := \frac{1}{m} \sum_{i=1}^m \log(1 + \exp(-y_i \theta^\top \tilde{x}_i)) + \frac{\lambda}{2} \|\theta\|^2. \quad (56)$$

其中

$$\ell(\theta) := \log(1 + \exp(-y \theta^\top \tilde{x})) \quad (57)$$

称为对数几率损失 (*logistic loss*).

Proof.

$$\begin{aligned} \theta^* &= \arg \min_{\theta} - \sum_{i=1}^m \log p(y_i | x_i, \theta) + \frac{\lambda}{2} \|\theta\|^2 \\ &= \arg \min_{\theta} - \sum_{i=1}^m \log \text{sigm } y_i \theta^\top \tilde{x}_i + \frac{\lambda}{2} \|\theta\|^2 \\ &= \arg \min_{\theta} \frac{1}{m} \sum_{i=1}^m \log(1 + \exp(-y_i \theta^\top \tilde{x}_i)) \\ &\quad + \frac{\lambda}{2} \|\theta\|^2. \end{aligned} \quad (58)$$

□

(3). 对损失函数 $\mathcal{L}(\theta)$ 进行梯度下降优化.

$$\begin{aligned} \frac{\partial \mathcal{L}}{\partial \theta} &= \frac{1}{m} \sum_{i=1}^m \frac{\exp(-y_i \theta^\top \tilde{x}_i)}{1 + \exp(-y_i \theta^\top \tilde{x}_i)} (-y_i x_i) + \lambda \theta \\ &= -\frac{1}{m} \sum_{i=1}^m \text{sigm}(-y_i \theta^\top \tilde{x}_i) y_i x_i + \lambda \theta. \end{aligned} \quad (59)$$

4 Softmax 回归

4.1 建模流程

Softmax 回归应对多分类问题, 它是对数几率回归向多分类问题的推广. 其建模方法包括如下三步 (参见第 1.1 节).

(1). 对 $p(y | x, \theta)$ 进行概率假设.

$$h(x)_k := p(y = k | x, \theta) := \frac{\exp \theta_k^\top \tilde{x}}{\sum_{c=1}^C \exp \theta_c^\top \tilde{x}} \in (0, 1). \quad (60)$$

对数几率回归假设 $p(y | x, \theta)$ 服从伯努利分布, Softmax 回归假设 $p(y | x, \theta)$ 服从如下分布

$$p(y | x, \theta) = \frac{\exp \theta_y^\top \tilde{x}}{\sum_{c=1}^C \exp \theta_c^\top \tilde{x}}. \quad (61)$$

令

$$\Theta := \begin{bmatrix} \theta_1^\top \\ \theta_2^\top \\ \vdots \\ \theta_C^\top \end{bmatrix} \in \mathbb{R}^{C \times (d+1)}, \quad (62)$$

假设函数可以写成矩阵的形式

$$p(y | x, \Theta) = \frac{(\exp \Theta \tilde{x})_y}{(\exp \Theta \tilde{x})^\top \mathbf{1}}. \quad (63)$$

(2). 对参数 θ 进行最大后验估计.

定理 15. 假设参数 θ 服从高斯先验, 对参数 θ 进行最大后验估计等价于优化如下损失函数

$$\mathcal{L}(\theta) = \frac{1}{m} \sum_{i=1}^m \left(-\theta_{y_i}^\top \tilde{x}_i + \log \sum_{c=1}^C \exp \theta_c^\top \tilde{x}_i \right) + \sum_{c=1}^C \frac{\lambda}{2} \|\theta_c\|^2. \quad (64)$$

其中

$$\ell(\theta) := -\theta_y^\top \tilde{x} + \log \sum_{c=1}^C \exp \theta_c^\top \tilde{x} \quad (65)$$

称为交叉熵损失 (*cross-entropy loss*).

Proof.

$$\begin{aligned} \theta^* &= \arg \min_{\theta} - \sum_{i=1}^m \log p(y_i | x_i, \theta) + \frac{\lambda}{2} \|\theta\|^2 \\ &= \arg \min_{\theta} - \sum_{i=1}^m \log \frac{\exp \theta_{y_i}^\top \tilde{x}_i}{\sum_{c=1}^C \exp \theta_c^\top \tilde{x}_i} + \sum_{c=1}^C \frac{\lambda}{2} \|\theta_c\|^2 \\ &= \arg \min_{\theta} - \sum_{i=1}^m \left(\theta_{y_i}^\top \tilde{x}_i - \log \sum_{c=1}^C \exp \theta_c^\top \tilde{x}_i \right) \\ &\quad + \sum_{c=1}^C \frac{\lambda}{2} \|\theta_c\|^2 \\ &= \arg \min_{\theta} \frac{1}{m} \sum_{i=1}^m \left(-\theta_{y_i}^\top \tilde{x}_i + \log \sum_{c=1}^C \exp \theta_c^\top \tilde{x}_i \right) \\ &\quad + \sum_{c=1}^C \frac{\lambda}{2} \|\theta_c\|^2. \end{aligned}$$

□

(3). 对损失函数 $\mathcal{L}(\theta)$ 进行梯度下降优化.

损失函数对应于类别 k 的参数 θ_k 的导数是

$$\begin{aligned} \frac{\partial \mathcal{L}}{\partial \theta_k} &= \frac{1}{m} \sum_{i=1}^m \left(-\mathbb{I}(y_i = k) x_i + \frac{\exp \theta_k^\top \tilde{x}_i}{\sum_{c=1}^C \exp \theta_c^\top \tilde{x}_i} x_i \right) + \lambda \theta_k \\ &= \frac{1}{m} \sum_{i=1}^m \left(\frac{\exp \theta_k^\top \tilde{x}_i}{\sum_{c=1}^C \exp \theta_c^\top \tilde{x}_i} - \mathbb{I}(y_i = k) \right) x_i + \lambda \theta_k. \end{aligned} \quad (66)$$

写成矩阵的形式是

$$\frac{\partial \mathcal{L}}{\partial \Theta} = \frac{1}{m} \sum_{i=1}^m \left(\frac{(\exp \Theta \tilde{\mathbf{x}}_i)^{y_i}}{(\exp \Theta \tilde{\mathbf{x}}_i)^{\top} \mathbf{1}} - \mathbf{e}_{y_i} \right) \mathbf{x}_i^{\top} + \lambda \Theta, \quad (67)$$

其中 $\mathbf{e}_k \in \mathbb{R}^C$ 的第 k 个元素是 1, 其余元素均为 0.

对比公式 20、49 和 67, 损失函数的梯度有相同的数学形式

$$\frac{\partial \mathcal{L}}{\partial \theta} = \frac{1}{m} \sum_{i=1}^m (h(\mathbf{x}_i) - y_i) \mathbf{x}_i + \lambda \theta. \quad (68)$$

区别在于假设函数 $h(\mathbf{x}_i)$ 的形式不同. 事实上, 所有的广义线性模型都有类似于公式 68 的更新准则.

4.2 交叉熵

定义由训练集观察得到的分布, 称为经验分布 (empirical distribution). 经验分布 p_i 对应于第 i 个样例, 定义 $p_{ic} := \mathbb{I}(y_i = c)$. 另一方面, $h(\mathbf{x}_i)$ 是由模型估计出的概率.

定理 16. 交叉熵损失旨在最小化经验分布 p_i 和学得分布 $h(\mathbf{x}_i)$ 之间的交叉熵. 这等价于最小化 p_i 和 $h(\mathbf{x}_i)$ 之间的 KL 散度, 迫使估计的分布 $h(\mathbf{x}_i)$ 近似目标分布 p_i .

Proof.

$$\begin{aligned} \theta^* &:= \arg \min_{\theta} -\frac{1}{m} \sum_{i=1}^m \log h(\mathbf{x}_i)_y + \frac{\lambda}{2} \|\theta\|^2 \\ &= \arg \min_{\theta} -\frac{1}{m} \sum_{i=1}^m \sum_{c=1}^C \mathbb{I}(y_i = c) \log h(\mathbf{x}_i)_c + \frac{\lambda}{2} \|\theta\|^2 \\ &= \arg \min_{\theta} -\frac{1}{m} \sum_{i=1}^m \sum_{c=1}^C p_{ic} \log h(\mathbf{x}_i)_c + \frac{\lambda}{2} \|\theta\|^2 \\ &= \arg \min_{\theta} \frac{1}{m} \sum_{i=1}^m \sum_{c=1}^C (p_{ic} \log p_{ic} - p_{ic} \log h(\mathbf{x}_i)_c) \\ &\quad + \frac{\lambda}{2} \|\theta\|^2 \\ &= \arg \min_{\theta} \frac{1}{m} \sum_{i=1}^m \sum_{c=1}^C p_{ic} \log \frac{p_{ic}}{h(\mathbf{x}_i)_c} + \frac{\lambda}{2} \|\theta\|^2 \\ &= \arg \min_{\theta} \frac{1}{m} \sum_{i=1}^m \text{KL}(p_i \parallel h(\mathbf{x}_i)) + \frac{\lambda}{2} \|\theta\|^2. \end{aligned} \quad (69)$$

□

5 朴素贝叶斯分类器

朴素贝叶斯分类器 (naive Bayes classifier) 也是一种概率方法, 但它是一种生成式模型. 在本节, 我们首先回顾生成式模型, 之后介绍朴素贝叶斯分类器的建模流程.

5.1 生成式模型

判别式模型和生成式模型各是什么? 判别式模型 (discriminant model) 直接对 $p(y | \mathbf{x})$ 进行建模, 生成式模型 (generative model) 先对联合分布 $p(\mathbf{x}, y) = p(\mathbf{x} | y)p(y)$ 进行建模, 然后再得到

$$p(y | \mathbf{x}) = \frac{p(\mathbf{x} | y)p(y)}{p(\mathbf{x})}. \quad (70)$$

其中, $p(y)$ 是类先验 (prior) 概率, 表达了样本空间中各类样本所占的比例. $p(\mathbf{x} | y)$ 称为似然 (likelihood). $p(\mathbf{x})$ 是用于归一化的证据 (evidence). 由于其和类标记无关, 该项不影响 $p(y | \mathbf{x})$ 的估计

$$\begin{aligned} h(\mathbf{x}) &:= \arg \max_c p(y = c | \mathbf{x}) \\ &= \arg \max_c p(\mathbf{x} | y = c)p(y = c) \\ &= \arg \max_c \log p(\mathbf{x} | y = c) + \log p(y = c) \end{aligned} \quad (71)$$

如何对类先验概率和似然进行估计? 根据大数定律, 当训练集包含充足的独立同分布样本时, $p(y)$ 可通过各类样本出现的频率来进行估计

$$p(y = c) := \frac{1}{m} \sum_{i=1}^m \mathbb{I}(y_i = c). \quad (72)$$

而对似然 $p(\mathbf{x} | y)$, 由于其涉及 \mathbf{x} 所有属性的联合概率, 如果基于有限训练样本直接估计联合概率, (1). 在计算上将会遭遇组合爆炸问题. (2). 在数据上将会遭遇样本稀疏问题, 很多样本取值在训练集中根本没有出现, 而“未被观测到”与“出现概率为零”通常是不同的. 直接按样本出现的频率来估计会有严重的困难, 属性数越多, 困难越严重.

判别式模型和生成式模型的优缺点? 优缺点对比如表 3 所示.

5.2 建模流程

(1). 对 $p(\mathbf{x} | y, \theta)$ 进行概率假设.

Table 3: 判别式模型和生成式模型的优缺点对比.

算法类型	算法性能	收敛速度	采样/生成服从潜在分布样例 (\mathbf{x}, y)	存在隐变量时进行学习
判别式模型	通常较高	通常较慢	不能	不能
生成式模型	通常较低	通常较快	能	能

生成式模型的主要困难在于, 类条件概率 $p(\mathbf{x} | y)$ 是所有属性的联合概率, 难以从有限的训练样本直接估计而得. 为避开这个障碍, 朴素贝叶斯分类器采用了属性条件独立性假设: 对已知类别, 假设所有属性相互独立. 也就是说, 假设每个属性独立地对分类结果发生影响

$$p(x_j, x_d | y) = p(x_j | y) p(x_d | y). \quad (73)$$

此外, 对连续属性, 进一步假设

$$p(x_j | y = c) \sim \mathcal{N}(\mu_{cj}, \sigma_{cj}^2). \quad (74)$$

因此, 朴素贝叶斯分类器的假设函数是

$$\begin{aligned} h(\mathbf{x}) &:= \arg \max_c \sum_{j=1}^d \log p(x_j | y = c) + \log p(y = c) \\ &= \arg \max_c \sum_{j=1}^d \log \frac{1}{\sqrt{2\pi\sigma_{cj}^2}} \exp\left(-\frac{(x_j - \mu_{cj})^2}{2\sigma_{cj}^2}\right) \\ &\quad + \log \sum_{i=1}^m \mathbb{I}(y_i = c) \\ &= \arg \min_c \sum_{j=1}^d \left(\frac{(x_j - \mu_{cj})^2}{2\sigma_{cj}^2} + \frac{1}{2} \log \sigma_{cj}^2 \right) \\ &\quad - \log \sum_{i=1}^m \mathbb{I}(y_i = c). \end{aligned} \quad (76)$$

(2). 对参数 θ 进行最大后验估计. 参数 θ 包括了第 c 类样本在第 j 个属性上的高斯分布的均值 μ_{cj} 和方差 σ_{cj}^2 .

定理 17. 假设参数 θ 服从不提供信息的先验, 对参数 θ 进行最大后验估计的结果是

$$\mu_{cj}^* := \frac{\sum_{i=1}^m \mathbb{I}(y_i = c) x_{ij}}{\sum_{i=1}^m \mathbb{I}(y_i = c)}; \quad (77)$$

$$\sigma_{cj}^{2*} := \frac{\sum_{i=1}^m \mathbb{I}(y_i = c) (x_{ij} - \mu_{cj}^*)^2}{\sum_{i=1}^m \mathbb{I}(y_i = c)}. \quad (78)$$

Proof. 代入公式 76 ,

$$\begin{aligned} \mu_{cj}^* &:= \arg \max_{\mu_{cj}} \sum_{i=1}^m \mathbb{I}(y_i = c) p(\mathbf{x}_i, y_i = c) \\ &= \arg \min_{\mu_{cj}} \sum_{i=1}^m \mathbb{I}(y_i = c) (x_j - \mu_{cj})^2. \end{aligned} \quad (79)$$

令 $\mathcal{L}(\mu_{cj}) := \sum_{i=1}^m \mathbb{I}(y_i = c) (x_j - \mu_{cj})^2$. 通过解

$$\frac{\partial \mathcal{L}}{\partial \mu_{cj}} = - \sum_{i=1}^m 2 \mathbb{I}(y_i = c) (x_j - \mu_{cj}) = 0 \quad (80)$$

可以得到 μ_{cj}^* .

$$\begin{aligned} \sigma_{cj}^2 &:= \arg \max_{\sigma_{cj}^2} \sum_{i=1}^m \mathbb{I}(y_i = c) p(\mathbf{x}_i, y_i = c) \\ &= \arg \min_{\sigma_{cj}^2} \sum_{i=1}^m \mathbb{I}(y_i = c) \left(\frac{(x_j - \mu_{cj})^2}{2\sigma_{cj}^2} + \frac{1}{2} \log \sigma_{cj}^2 \right). \end{aligned}$$

令 $\mathcal{L}(\sigma_{cj}^2) := \sum_{i=1}^m \mathbb{I}(y_i = c) \left(\frac{(x_j - \mu_{cj})^2}{2\sigma_{cj}^2} + \frac{1}{2} \log \sigma_{cj}^2 \right)$. 通过解

$$\frac{\partial \mathcal{L}}{\partial \sigma_{cj}^2} = \sum_{i=1}^m \mathbb{I}(y_i = c) \left(-\frac{(x_j - \mu_{cj})^2}{2\sigma_{cj}^4} + \frac{1}{2} \frac{1}{\sigma_{cj}^2} \right) = 0 \quad (81)$$

可以得到 σ_{cj}^{2*} . \square

5.3 离散属性的参数估计

朴素贝叶斯分类器可以很容易地处理离散属性. $p(x_j | y)$ 可估计为

$$p(x_j = a | y = c) := \frac{\sum_{i=1}^m \mathbb{I}(x_j = a \wedge y_i = c)}{\sum_{i=1}^m \mathbb{I}(y_i = c)}. \quad (82)$$

然而, 若某个属性值在训练集中没有与某个类同时出现过, 则根据公式 82 估计得到 0. 代入公式 75 得到 $-\infty$. 因此, 无论该样本的其他属性是什么, 分类结果都不会是 $y = c$, 这显然不太合理.

为了避免其他属性携带的信息被训练集中未出现的属性值“抹去”, 在估计概率值时通常要进行平滑

(smoothing), 常用拉普拉斯修正 (Laplacian correction). 具体的说, 令 K 表示训练集 D 中可能的类别数, n_j 表示第 j 个属性可能的取值数, 则概率估计修正为

$$p(y = c) := \frac{\sum_{i=1}^m \mathbb{I}(y_i = c) + 1}{m + K}; \quad (83)$$

$$p(x_j = a | y = c) := \frac{\sum_{i=1}^m \mathbb{I}(x_j = a \wedge y_i = c) + 1}{\sum_{i=1}^m \mathbb{I}(y_i = c) + n_j}. \quad (84)$$

拉普拉斯修正实际上假设了属性值与类别均匀分布, 这是在朴素贝叶斯学习中额外引入的关于数据的先验. 在训练集变大时, 修正过程所引入的先验的影响也会逐渐变得可忽略, 使得估值渐趋向于实际概率值.

在现实任务中朴素贝叶斯有多种实现方式. 例如, 若任务对预测速度要求较高, 则对给定训练集, 可将朴素贝叶斯分类器涉及的所有概率估值事先计算好存储起来, 这样在进行预测时只需查表即可进行判别. 若任务数据更替频繁, 则可采用懒惰学习方式, 先不进行任何训练, 待收到预测请求时再根据当前数据集进行概率估值. 若数据不断增加, 则可在现有估值基础上, 仅对新增样本的属性值所涉及的概率估值进行计数修正即可实现增量学习.

定义 3 (懒惰学习 (lazy learning)). 这类学习技术在训练阶段仅仅是把样本保存起来, 训练时间开销是 0, 待收到测试样本后再进行处理. 相应的, 那些在训练阶段就对样本进行学习处理的方法称为急切学习 (eager learning).

定义 4 (增量学习 (incremental learning)). 在学得模型后, 再接收到训练样例时, 仅需根据新样例对模型进行更新, 不必重新训练整个模型, 并且先前学得的有效信息不会被“冲掉”.

5.4 朴素贝叶斯分类器的推广

朴素贝叶斯分类器采用了属性条件独立性假设, 但在现实任务中这个假设往往很难成立. 于是, 人们尝试对属性条件独立性假设进行一定程度的放松, 适当考虑一部分属性间的相互依赖关系, 这样既不需要进行完全联合概率计算, 又不至于彻底忽略了比较强的属性依赖关系, 由此产生一类半朴素贝叶斯分类器 (semi-naive Bayes classifiers) 的学习方法.

独依赖估计 (one-dependent estimator, ODE) 是最常用的一种策略, 其假设每个属性在类别之外最多依赖

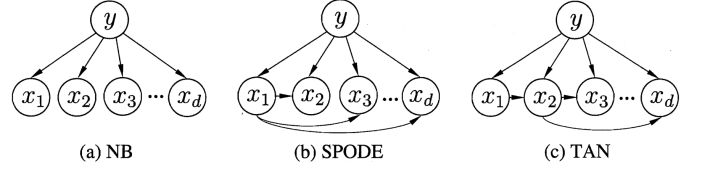


Figure 6: 朴素贝叶斯 (NB) 和两种半朴素贝叶斯分类器 (SPODE、TAN) 所考虑的属性依赖关系. 本图源于 [17].

于一个其他属性 (称为父属性). 问题的关键在于如何确定每个属性的父属性. SPODE (super-parent ODE) 假设所有属性都依赖于同一个属性, 称为超父 (super-parent). TAN (tree augmented naive Bayes) [4] 以属性节点构建完全图, 任意两结点之间边的权重设为这两个属性之间的条件互信息 $I(x_j, x_k | y)$. 之后构建此图的最大带权生成树, 挑选根变量, 将边置为有向, 以将属性间依赖关系约简为树形结构. 最后加入类别结点 y , 增加从 y 到每个属性的有向边. TAN 通过条件互信息刻画两属性的条件相关性, 最终保留了强相关属性之间的依赖性. AODE (averaged ODE) [13] 尝试将每个属性作为超父来构建 SPODE, 之后将那些具有足够训练数据支撑的 SPODE 集成作为最终结果. AODE 的训练过程也是“计数”, 因此具有朴素贝叶斯分类器无需模型选择、可预计算节省预测时间、也能懒惰学习、并且易于实现增量学习.

能否通过考虑属性间高阶依赖进一步提升泛化性能? 相比 ODE, k DE 考虑最多 k 个父属性. 随着依赖的属性个数 k 的增加, 准确进行概率估计所需的训练样本数量将以指数级增加. 因此, 若训练数据非常充分, 泛化性能有可能提升. 但在有限样本条件下, 则又陷入高阶联合概率的泥沼.

更进一步, 贝叶斯网 (Bayesian network), 也成为信念网 (belief network), 能表示任意属性间的依赖性. 贝叶斯网是一种概率图模型, 借助有向无环图刻画属性间的依赖关系.

事实上, 虽然朴素贝叶斯的属性条件独立假设在现实应用中往往很难成立, 但在很多情形下都能获得相当好的性能 [2, 8]. 一种解释是对分类任务来说, 只需各类别的条件概率排序正确, 无须精准概率值即可导致正确分类结果 [2]. 另一种解释是, 若属性间依赖对所有类别影响相同, 或依赖关系能相互抵消, 则属性条件独立性假设在降低计算开销的同时不会对性能产生负面

影响 [15]。朴素贝叶斯分类器在信息检索领域尤为常用 [6]。

6 快问快答

随机梯度下降和标准梯度下降的优缺点各是什么？

- 参数更新速度. 标准梯度下降需要遍历整个训练集才能计算出梯度, 更新较慢. 随机梯度下降只需要一个训练样例即可计算出梯度, 更新较快.
- 冗余计算. 当训练集样本存在冗余时, 随机梯度下降能避免在相似样例上计算梯度的冗余.
- 梯度中的随机因素/噪声. 标准梯度下降计算得到的梯度没有随机因素, 一旦陷入局部极小将无法跳出. 随机梯度下降计算得到的梯度有随机因素, 有机会跳出局部极小继续优化.

实际应用时, 常采用随机梯度下降和标准梯度下降的折中, 即使用一部分样例进行小批量梯度下降. 此外, 相比随机梯度下降, 小批量梯度下降还可以更好利用矩阵的向量化计算的优势.

梯度下降和牛顿法的优缺点各是什么？

- 导数阶数. 梯度下降只需要计算一阶导数, 而牛顿法需要计算二阶导数. 一阶导数提供了方向信息 (下降最快的方向), 二阶导数还提供了函数的形状信息.
- 计算和存储开销. 牛顿法在参数更新时需要计算 Hessian 矩阵的逆, 计算和存储开销比梯度下降更高.
- 学习率. 梯度下降对学习率很敏感, 而标准的牛顿法不需要设置学习率.
- 收敛速度. 牛顿法的收敛速度比梯度下降更快.
- 牛顿法不适合小批量或随机样本.

实际应用时, 有许多拟牛顿法旨在以较低的计算和存储开销近似 Hessian 矩阵.

线性回归的损失函数及梯度推导.

答案见上文.

为什么要使用正则化, ℓ_1 和 ℓ_2 正则化各自对应什么分布, 各有什么作用?

答案见上文.

对数几率回归的损失函数及梯度推导.

答案见上文.

对数几率回归和支持向量机的相同点和不同点?

相同点:

- 对数几率回归和支持向量机的主问题均为参数模型.
- 损失函数均为经验风险 + 结构风险 (正则化项) 的形式, 其经验风险均是 0/1 损失函数的替代损失, 结构风险都可选用 ℓ_1 、 ℓ_2 正则化项.
- 它们均为凸优化问题, 能有收敛到全局最优的保证. 对数几率回归和支持向量机的主问题均常采用基于梯度的方法进行优化.
- 均为二分类线性学习算法, 均可扩展到处理非线性、多分类问题.

不同点:

- 支持向量机的对偶问题是非参数模型, 需要利用训练集中的支持向量.
- 经验风险的函数形式不同 (对数几率损失和合页损失).
- 支持向量机的对偶问题采用二次规划进行优化, 或使用专门设计的 SMO 优化算法.
- 支持向量机不是概率方法, 不需要对 $p(y | \mathbf{x}, \theta)$ 进行概率假设. 对数几率回归是判别式模型, 而支持向量机直接得到判别函数.
- 对数几率回归是广义线性模型, 而支持向量机不是广义线性模型.

线性分类器如何扩展为非线性分类器?

答案见上文.

判别式模型和生成式模型各是什么, 各自优缺点是什么, 常见算法中哪些是判别式模型, 哪些是生成式模型?

答案见上文.

贝叶斯定理各项的含义?

答案见上文.

朴素贝叶斯为什么叫“朴素”贝叶斯?

为了避免从有限的训练样本直接估计 $p(\mathbf{x} | y)$ 的障碍, 朴素贝叶斯做出了属性条件独立假设, 该假设在现

实应用中往往很难成立.

References

- [1] P. L. Combettes and V. R. Wajs. Signal recovery by proximal forward-backward splitting. *Multiscale Modeling & Simulation*, 4(4):1168–1200, 2005. 5
- [2] P. M. Domingos and M. J. Pazzani. On the optimality of the simple bayesian classifier under zero-one loss. *Machine Learning*, 29(2-3):103–130, 1997. 12
- [3] B. Efron. Bayesians, frequentists, and scientists. *Journal of the American Statistical Association*, 100(469):1–5, 2005. 1
- [4] N. Friedman, D. Geiger, and M. Goldszmidt. Bayesian network classifiers. *Machine Learning*, 29(2-3):131–163, 1997. 12
- [5] P. J. Huber. Robust estimation of a location parameter. *Annals of Statistics*, 53(1):492–518, 1964. 6
- [6] D. D. Lewis. Naive (bayes) at forty: The independence assumption in information retrieval. In *Proceedings of the 10th European Conference on Machine Learning (ECML)*, pages 4–15, 1998. 13
- [7] K. P. Murphy. *Machine Learning: A Probabilistic Perspective*. MIT Press, 2012. 5, 6
- [8] A. Y. Ng and M. I. Jordan. On discriminative vs. generative classifiers: A comparison of logistic regression and naive bayes. In *Advances in Neural Information Processing Systems 14 (NIPS)*, pages 841–848, 2001. 12
- [9] F. J. Samaniegos. *A Comparison of the Bayesian and Frequentist Approaches to Estimation*. Springer Science & Business Media, 2010. 1
- [10] R. Tibshirani. Regression shrinkage and selection via the LASSO. *Journal of the Royal Statistical Society. Series B (Methodological)*, pages 267–288, 1996. 4
- [11] R. Tibshirani, M. Saunders, S. Rosset, J. Zhu, and K. Knight. Sparsity and smoothness via the fused lasso. *Journal of the Royal Statistical Society: Series B (Statistical Methodology)*, 67(1):91–108, 2005. 5
- [12] A. N. Tikhonov and V. I. Arsenin. *Solutions of Ill-posed Problems*. Winston, 1977. 4
- [13] G. I. Webb, J. R. Boughton, and Z. Wang. Not so naive bayes: Aggregating one-dependence estimators. *Machine Learning*, 58(1):5–24, 2005. 12
- [14] M. Yuan and Y. Lin. Model selection and estimation in regression with grouped variables. *Journal of the Royal Statistical Society: Series B (Statistical Methodology)*, 68(1):49–67, 2006. 5
- [15] H. Zhang. The optimality of naive bayes. In *Proceedings of the Seventeenth International Florida Artificial Intelligence Research Society Conference (FLAIRS)*, pages 562–567, 2004. 13
- [16] H. Zou and T. Hastie. Regularization and variable selection via the elastic net. *Journal of the Royal Statistical Society: Series B (Statistical Methodology)*, 67(2):301–320, 2005. 5
- [17] 周志华. 机器学习. 清华大学出版社, 2016. 5, 7, 12