

四、5、贝叶斯估计

2016年3月5日 星期六

下午12:53

简介

贝叶斯估计认为，即使训练集已知，待估计的参数仍然是未知的，但是待估计的参数服从某种特定的分布，因此可以假设待估计参数服从一个先验的分布，训练数据作为观测数据，见过观测数据之后，对待估计参数有一个更准确的估计（把先验概率转化成后验概率密度），这个估计就是我们待求的参数。这个方法我们称作贝叶斯估计。

我的理解，之所以叫做贝叶斯估计，是因为 1、参数估计的思想借鉴贝叶斯学派 2、参数估计的过程是使用贝叶斯公式，通过逐步增加观测数据，把参数 θ 的先验概率变成后验概率，最终产出一个更准确的参数 θ 的概率分布。

算法原理

整体思路

第一步、贝叶斯估计：

使用贝叶斯公式，通过逐步增加观测数据 $x \in D$ ，把参数 θ 的先验概率 $P(\theta)$ 变成后验概率 $P(\theta | D)$ ，从而产出一个更准确的参数 θ 的概率分布。这个更准确的概率分布其实就是后验概率 $P(\theta | D)$ ，得到参数 θ 的概率分布，贝叶斯估计就已经完成了。

第二步、计算类条件概率：

在得到参数 θ 更准确的概率分布基础上，往往要结合实际问题的，用到实际问题中。这里以贝叶斯分类器为例，我们要做一个贝叶斯分类器，最终目标是要通过参数 θ 估计的是类条件概率：对于第 i 类： $P(x | w_i, \theta_i)$ ，即 $P(x | \theta_i)$ 。一

一般省去类下标写成 $P(x|\theta)$ 。

这里可以回顾一下，类条件概率的作用是带入贝叶斯分类器的贝叶斯公式中，计算出后验概率。后验概率是贝叶斯分类器分类的依据

由于这里我们得到的是 θ 的概率分布，（而不仅仅是最大似然估计中，得到参数的一个固定值），所以我们计算类条件概率时可以充分利用“概率分布”的强大之处。我们对 θ 在各种取值情况下产出的类条件概率做积分。

$$P(x|D) = \int P(x|\theta) P(\theta|D) d\theta.$$

注意这里类条件概率写成 $P(x|D)$ 而非 $P(x|\theta)$ 是因为 θ 是一个概率分布，不是固定值，之前的“在参数 θ 情况下， x 的概率”写成“在数据集 D 情况下， θ 在各种取值下的 x 的概率的积分”。其中 $P(\theta|D)$ 已经估计出， $P(x|\theta)$ 后面会简化，至此可得出类条件概率 $P(x|D)$

求解目标

第一步： $P(\theta)$ 的后验概率 $P(\theta|D)$

第二步：类条件概率 $P(x|D)$ （通常表示为 $P(x|\theta)$ ，只不过贝叶斯估计比较特殊）

已知与假设

$P(\theta)$

我们要估计 $P(\theta)$ 的概率（其实是概率分布），需要先假定 $P(\theta)$ 服从概率分布的形式，我们假定他服从高斯分布。

$$P(\theta) \sim N(\mu_0 | \sigma_0^2)$$

$P(x|\theta)$

假设 θ 中协方差 Σ 是已知的，只有 μ 是未知的。

而且 $P(x|\mu) \sim N(\mu | \Sigma)$ 服从正态分布

这里两个 μ 是一样的，是正好一致吗，可能是，没太理解，详

单变量情况

我们先讨论单变量情况，（变量 x 等是一维的，而不是多维特征），多变量情况其实是单变量基础上使用向量、矩阵乘法得到的详见模式分类P77。

求解过程（尤其在特殊情况下）

第一步： $P(\theta|D)$

贝叶斯公式求解过程

由于 θ 中只有 μ 是未知的，表示为 $P(\mu|D)$ ，根据贝叶斯公式

$$P(\mu|D) = P(D|\mu) * P(\mu) / P(D)$$

其中 $P(D|\mu) = \prod P(x_k|\mu)$ 对于 D 中的 n 个样本 x_k ， \prod 是连乘符号

$P(x_i|\mu)$ 服从高斯分布

$P(\mu)$ 服从高斯分布

$P(D)$ 看做常量，记作 $1/a$

进而

$$\begin{aligned} p(\mu|D) &= \alpha \prod_{k=1}^n \overbrace{\frac{1}{\sqrt{2\pi}\sigma} \exp\left[-\frac{1}{2}\left(\frac{x_k - \mu}{\sigma}\right)^2\right]}^{p(x_k|\mu)} \overbrace{\frac{1}{\sqrt{2\pi}\sigma_0} \exp\left[-\frac{1}{2}\left(\frac{\mu - \mu_0}{\sigma_0}\right)^2\right]}^{p(\mu)} \\ &= \alpha' \exp\left[-\frac{1}{2}\left(\sum_{k=1}^n \left(\frac{\mu - x_k}{\sigma}\right)^2 + \left(\frac{\mu - \mu_0}{\sigma_0}\right)^2\right)\right] \\ &= \alpha'' \exp\left[-\frac{1}{2}\left[\left(\frac{n}{\sigma^2} + \frac{1}{\sigma_0^2}\right)\mu^2 - 2\left(\frac{1}{\sigma^2}\sum_{k=1}^n x_k + \frac{\mu_0}{\sigma_0^2}\right)\mu\right]\right] \end{aligned}$$

两个高斯分布经过化简成为了一个高斯分布， $P(\mu|D)$ 也服从高斯分布。 $P(\mu|D)$ 的均值 μ_n 和方差 σ_n 可以用之前那两个高斯分布的均值方差表示为

$$\frac{1}{\sigma_n^2} = \frac{n}{\sigma^2} + \frac{1}{\sigma_0^2}$$

$$\mu_n = \left(\frac{n\sigma_0^2}{n\sigma_0^2 + \sigma^2} \right) \hat{\mu}_n + \frac{\sigma^2}{n\sigma_0^2 + \sigma^2} \mu_0$$

其中 μ_n 是数据集样本均值

$$\hat{\mu}_n = \frac{1}{n} \sum_{k=1}^n x_k$$

分析：

从结果中，我们能看出来对于待估计参数 θ 先验知识和样本观测是如何结合的。

μ_n 是使得 θ 概率最大的参数估计， σ_n 是经过多次观测之后，仍然留有的对估计的不确定程度。可以看出

A σ_n 随样本 n 增加是单减的。

B μ_n 随样本增加趋近样本均值

C 随着样本增加 $P(\mu | D)$ 越来越尖，这个过程称为贝叶斯学习过程。

D 特殊情况： $\sigma_0 = 0$ ，说明先验十分置信，任何样本都无法改变后验概率。

E 特殊情况： $\sigma_0 \gg \sigma$ ，说明先验十分不置信，直接把样本 μ 当做参数的后验的均值即可。

第二步： $P(x | D)$

得到参数估计 $P(\mu | D)$ 之后，我们要求类条件概率 $P(x | D)$ （其他算法中表示为 $P(\mu | \theta)$ ）

利用贝叶斯估计的优势：求得的参数 θ 是一个概率分布，而非特定值。求解 $P(x | D)$ 时要结合概率分布，所以：

$$P(x | D) = \int P(x | \theta) P(\theta | D) d_\theta$$

由于参数 θ 只有 μ 未知

$$P(x | D) = \int P(x | \mu) P(\mu | D) d_\mu$$

$P(x | \mu)$ 服从高斯分布

$P(\mu | D)$ 我们已经求得，而且也服从高斯分布

$$\begin{aligned}
 p(x|\mathcal{D}) &= \int p(x|\mu) p(\mu|\mathcal{D}) d\mu \\
 &= \int \frac{1}{\sqrt{2\pi}\sigma} \exp\left[-\frac{1}{2}\left(\frac{x-\mu}{\sigma}\right)^2\right] \frac{1}{\sqrt{2\pi}\sigma_n} \exp\left[-\frac{1}{2}\left(\frac{\mu-\mu_n}{\sigma_n}\right)^2\right] d\mu \\
 &= \frac{1}{2\pi\sigma\sigma_n} \exp\left[-\frac{1}{2}\frac{(x-\mu_n)^2}{\sigma^2+\sigma_n^2}\right] f(\sigma, \sigma_n)
 \end{aligned}$$

其中

$$f(\sigma, \sigma_n) = \int \exp\left[-\frac{1}{2}\frac{\sigma^2+\sigma_n^2}{\sigma^2\sigma_n^2}\left(\mu - \frac{\sigma_n^2x + \sigma^2\mu_n}{\sigma^2 + \sigma_n^2}\right)^2\right] d\mu$$

所以 $P(x|\mathcal{D})$ 也服从高斯分布

$$p(x|\mathcal{D}) \sim N(\mu_n, \sigma^2 + \sigma_n^2)$$

至此得到 $P(x|\mathcal{D})$

最终结果上，贝叶斯估计得到的是类条件概率的概率分布，最大似然估计只得到均值和方差的值