

# 贝叶斯后验的快速计算方法

闵素芹<sup>1</sup>, 李群<sup>2</sup>

(1.中国传媒大学 理学院,北京 100024;2.中国社会科学院 数量经济与技术经济研究所,北京 100732)

**摘要:**贝叶斯推断的主要障碍是后验分布积分计算,INLA方法通过对边缘后验密度的精确近似克服了计算时间的问题,而且得到的计算结果与MCMC的精度几乎无异。文章首先对INLA方法的计算进行介绍,然后通过综述INLA与MCMC的比较研究阐述了INLA方法应用的优越性,并给出一个运用INLA进行参数估计的应用实例。

**关键词:**贝叶斯推断;后验边缘;INLA;MCMC

**中图分类号:** O212.8 **文献标识码:** A **文章编号:** 1002-6487(2017)22-0084-03

## 0 引言

贝叶斯推断需要计算后验分布的积分,而后验分布往往是复杂的、高维的、非标准形式的分布,积分难以计算。MCMC方法是常用的一种行之有效的贝叶斯计算方法,该方法渐进精确,但是存在运行速度相对较慢的问题,特别是在大数据背景下很难适应新的模型的要求。

INLA是由Rue等(2009)<sup>[1]</sup>针对结构可加回归模型提出的一种近似贝叶斯推断的快速方法,它基于近似数值积分计算贝叶斯后验边缘。该方法的估计精度与MCMC方法类似,但在计算时间上远远快于MCMC方法,对于INLA需要几秒到几分钟的问题,MCMC通常需要几分钟到几小时。而且R-INLA程序包目前已比较成熟,应用方便。

INLA方法提出后由于其运算速度上突出的优越性得到了广泛关注,应用领域包括广义线性模型、平滑样条模型、半参数回归、空间和时空模型、地统计模型等。例如,Martino等(2010)<sup>[2]</sup>运用INLA对金融时间序列的随机波动模型(SV模型)进行推断,基于1217天SP500指数日收盘数据和1292天微软收盘价数据进行实际数据分析。Fong等(2010)<sup>[3]</sup>运用INLA方法对广义线性混合模型(GLMM)的贝叶斯推断进行了研究。Yu和Rue(2011)<sup>[4]</sup>运用INLA方法对分位数 $\tau=0.25,0.50$ 和 $0.75$ 分别进行参数估计,针对慕尼黑房屋租赁数据,因变量为每平米的租金,模型包括13个自变量的线性效应、面积和建筑年份的非线性效应、以及房子位置的空间效应。Natário等(2014)<sup>[5]</sup>利用INLA方法运用时空分层模型对葡萄牙森林火灾数据进行分析。Gómez-Rubio(2015)<sup>[6]</sup>利用INLA方法对空间计量模型参数进行推断,分析波士顿房产数据。

近几年,INLA方法在应用领域突出出优越的快速计

算能力。本文对其计算方法、优点进行阐述,并给出一个应用实例,旨在为INLA方法的进一步研究和应用提供借鉴。

## 1 INLA方法

INLA算法由Rue等(2009)<sup>[1]</sup>提出,该算法应用于一类应用广泛的结构可加回归模型的贝叶斯推断。模型的因变量假定服从指数分布族,均值 $\mu_i$ 通过一个连接函数 $g(\cdot)$ 连接到一个结构可加的预测变量 $\eta_i$ ,即: $g(\mu_i)=\eta_i$ 。

$$\eta_i = \alpha + \sum_{j=1}^{n_f} f^{(j)}(u_{ji}) + \sum_{k=1}^{n_B} \beta_k z_{ki} + \varepsilon_i \quad (1)$$

其中, $\{f^{(j)}(\cdot)\}$ 是自变量 $u$ 的函数,可选用不同的形式,如:自回归模型、季节效应、空间效应等, $\{\beta_k\}$ 代表自变量 $z$ 的线性效应, $\varepsilon_i$ 是非结构项。

结构可加回归模型是一类应用非常广泛的模型,包括广义线性模型、平滑样条模型、状态空间模型、半参数回归、空间和时空模型、对数高斯Cox过程和地统计模型等。

假定模型中的 $\alpha, \{f^{(j)}\}, \{\beta_k\}, \{\varepsilon_i\}$ 具有高斯先验。令向量 $x$ 表示模型中 $\eta_i, \alpha, \{f^{(j)}\}, \{\beta_k\}$ 这 $n$ 个变量, $\pi(\cdot|\cdot)$ 表示参数的条件密度,假定 $\pi(x|\theta_1)$ 是零均值、精度矩阵 $Q(\theta_1)$ 的高斯密度, $n_d$ 个观测变量 $y=\{y_i:i \in I\}$ 的分布用 $\pi(y|x, \theta_2)$ 表示,假定 $y=\{y_i:i \in I\}$ 在给定 $x, \theta_2$ 时条件独立,令 $\theta=(\theta_1^T, \theta_2^T)$ , $\theta$ 是 $m$ 维向量。

主要目标是计算后验边缘密度 $\pi(x_i|y)$ ,  $\pi(\theta|y)$ 和 $\pi(\theta_j|y)$ 。

边缘后验密度记为:

$$\pi(x_i|y) = \int \pi(x_i|\theta, y) \pi(\theta|y) d\theta$$

**基金项目:**北京高等学校“青年英才计划”项目(YETP0611);中国传媒大学优秀中青年教师培养工程(YXJS201330)

**作者简介:**闵素芹(1978—),女,山东青州人,博士,副教授,研究方向:分层模型、空间统计理论。

李群(1961—),男,山东聊城人,教授,博士生导师,研究方向:计量经济模型。

$$\pi(\theta_j|y) = \int \pi(\theta|y) d\theta_{-j}$$

INLA 方法构建一个嵌套近似:

$$\tilde{\pi}(x_i|y) = \int \tilde{\pi}(x_i|\theta, y) \tilde{\pi}(\theta|y) d\theta \quad (2)$$

$$\tilde{\pi}(\theta_j|y) = \int \tilde{\pi}(\theta|y) d\theta_{-j}$$

通过近似  $\pi(\theta|y)$  和  $\pi(x_i|\theta, y)$ , 运用数值积分来计算  $\pi(x_i|y)$  的近似。INLA 方法近似边缘后验分为 3 步来完成:

(1) 基于如下拉普拉斯近似计算  $\theta$  的边缘后验:

$$\tilde{\pi}(\theta|y) \propto \frac{\pi(x, \theta, y)}{\tilde{\pi}_G(x|\theta, y)} \Bigg|_{x=x^*(\theta)} \quad (3)$$

其中,  $\tilde{\pi}_G(x|\theta, y)$  是  $x$  的完全条件高斯近似,  $x^*(\theta)$  是给定  $\theta$  时  $x$  的众数。

(2) 计算  $x_i|\theta, y$  密度的拉普拉斯近似:

$$\tilde{\pi}(x_i|\theta, y) = N\{x_i; \mu_i(\theta), \sigma_i^2(\theta)\}$$

其中,  $\mu(\theta)$  是高斯近似的均值,  $\sigma^2(\theta)$  是相应的边缘方差向量。

(3) 根据式(2)对  $\theta$  进行数值积分, 通过下式将前两步结合起来:

$$\tilde{\pi}(x_i|y) = \sum_k \tilde{\pi}(x_i|\theta_k, y) \tilde{\pi}(\theta_k|y) \Delta_k$$

INLA 方法是一种集成的嵌套拉普拉斯近似, 联合解析近似和数值积分得到近似后验边缘, 可以用来计算后验均值、方差和分位数等。该方法实现了快速计算, 代价是解析近似在计算后验概率时可能会有误差。另外, INLA 计算的是边缘后验而非联合后验。但模拟研究均表明, 在提高计算速度的同时, INLA 方法可以保证较高的估计精度。

## 2 INLA 与 MCMC 的比较

由于 MCMC 方法是目前贝叶斯推断中应用最为广泛的一种算法, INLA 通过与其比较, 可以直观的阐明 INLA 方法的优越性。

在对模型进行贝叶斯推断时需要计算后验分布的积分。MCMC 方法的研究为推广贝叶斯推断方法的应用做出了贡献, 使得之前不能实施计算的一些统计方法变得比较容易。求普通方法无法得到的后验分布密度常常运用高维积分运算, MCMC 方法可以得到一条或几条收敛的马尔可夫链, 该马尔可夫链的极限分布就是所需的后验分布。MCMC 方法的基本思想是通过建立一个平稳分布为  $p(\theta|x^n)$  的马尔可夫链来得到  $p(\theta|x^n)$  的样本, 基于这些样本就可以作各种统计推断<sup>[7]</sup>。

然而, 结构可加回归模型运用 MCMC 方法时暴露出两个缺点, 一是  $x$  中的各分量之间强相关, 二是  $\theta$  和  $x$  也相关, 尤其当  $n$  很大时。而 INLA 方法在收敛问题和运算速度上得到很大改善, 该方法提出后, 很多学者通过模拟数据或实际数据对其运行时间和估计精度与 MCMC 方法

进行了比较研究。

Paul M 等(2010)<sup>[8]</sup>基于端粒酶数据, 运用广义双随机效应 meta 分析模型, 对 INLA (运用 R 软件的 INLA 程序包) 和 MCMC (运用 SAS 软件的 PROC NLMIXED) 进行比较, INLA 运行时间小于 0.2 秒, 而 MCMC 通过 10000 次后验抽样 (5000 次预迭代), 运行时间 10 分钟。通过进一步模拟研究比较两者的运行时间和失败次数 (即: 产生诸如优化数值计算不收敛所导致的不可靠结果)。每个 meta 分析选用 25 个研究, 共模拟 72 种不同的情况, 每种情况由 1000 次 meta 分析组成。所有分析均在 Intel(R) Core(TM) 2 Duo T7200 processor 2.00 GHz 进行, INLA 方法平均花费只有 3.8 分钟 (2.2 ~ 6.0 分钟), 而 MCMC 方法平均花费 42.2 分钟 (23.1 ~ 74.7 分钟)。对于不收敛率, INLA 方法非常稳定, 仅失败两次; 而 MCMC 方法收敛问题与 meta 分析的设计相关, 总体来看, 不收敛率为 1.6% (72000 次分析中失败 409 次), 并且 72000 次中有 11037 次的相关系数标准误差非常大。另外, 两种方法得到的灵敏度 (即真正率) 和特指度 (即真负率) 的偏倚和 MSE 表现非常接近。综合来看, 基于 MCMC 抽样的贝叶斯方法往往需要很长的运行时间, INLA 方法快速, 并且估计更稳健、具有更高的精度和更小的偏倚。

Yu 和 Rue (2011)<sup>[9]</sup>针对包含非线性项和随机效应项的可加混合分位数回归模型对 INLA 和 MCMC 进行比较, 模拟生成了两个数据集, 样本量均为 400, 重复 200 次, 对每个数据集估计 0.10, 0.25, 0.50, 0.75, 0.90 分位数。MCMC 方法基于 40000 次迭代 (15000 次预迭代)。采用平均绝对偏差 (MADE) 对估计进行评估, INLA 和 MCMC 的模拟研究结果几乎无异。并且, 当估计极端的分位数时 MCMC 逊于 INLA。为了比较运行时间, 模拟了样本量 400, 800, 1500, 3000 和 6000 的情况。INLA 运行速度明显更快, 各样本量下的计算时间分别为 1.01 秒、1.97 秒、3.84 秒、8.25 秒、18.09 秒; 而 MCMC 方法相应为 22.64 秒、46.82 秒、87.89 秒、172.29 秒、338.99 秒, MCMC 计算速度随着样本量的增加而严重变缓。其中, INLA 在 single-processor 2.13-GHz laptop 上运行, 运用 R-INLA; MCMC 在 quad-core Intel Pentium 2.84 GHz CPUs 上运行, 运用 Fortran 程序。

Grilli L 等 (2015)<sup>[9]</sup>针对二元 logit 混合模型, 通过模拟研究对 INLA 和 MCMC (Gibbs 抽样器) 的贝叶斯推断进行比较。对不同的先验设定, INLA 方法总是收敛的更快, 计算迅速。当组数 (层-2 单位数) 较大时 (大于 70) 所有的先验设定均能得到较好的估计结果, 先验的选择会影响 INLA 和 MCMC 的估计。模拟不同的组数和每组观测数时组方差估计的偏倚, 最坏的情况 (10 组, 每组 10 个观测) 下, INLA 的偏倚仅增加 10%, 但 MCMC 增加 20.8%。回归系数点估计的偏倚, 两种方法的结果类似。综合来看, INLA 比 MCMC 更精确, 计算时间上更优越, 例如: 组数 100, 每组观测数 50 时, INLA 耗时 11 秒, 而 MCMC 却需要 534 秒。

Schrödle 等 (2011)<sup>[10]</sup>基于牛病毒性腹泻数据的时空模型, 对所有参数的后验分布 MCMC 方法得到的直方图与

INLA 得到的近似曲线是一致的。MCMC 抽样所得的 DIC 值与由 INLA 计算的 DIC 值非常接近。INLA 计算时间为 233.23 秒,而 MCMC(3030000 次迭代,30000 次预迭代)速度为每秒进行 246 次迭代,即:花费时间远远高于 INLA。

Taylor(2012)<sup>[11]</sup>针对空间连续高斯过程,通过模拟研究比较 INLA 与 MCMC。MCMC 运用 M-H 算法,进行 100000 次迭代(10000 次预迭代)。通过不同的参数设置模拟了 18 种情况,综合来看,INLA 计算时间降低了 30%~50%,两种方法 MSE 无显著差异。

### 3 应用实例

下面通过一个简单的例子来说明如何运用 INLA 方法进行参数估计。以微博评论数影响因素研究数据为例,因变量为  $y$  (评论数),微博的评论数是计数数据,数据非负、分布尖峰、严重偏斜、方差明显大于均值、过度离散,因此适用广义线性模型中的负二项回归。自变量为每条微博具有的特点,包括  $x_1$  (时新性)、 $x_2$  (时间性)、 $x_3$  (重要性)、 $x_4$  (趣味性)、 $x_5$  (延续性)、 $x_6$  (显著性)、 $x_7$  (个体性)、 $x_8$  (正面性)。 $x_1$  和  $x_2$  是多分类数据,转换为虚拟变量来处理; $x_3-x_8$  为二分类数据。

运用 R 软件 INLA 程序包中的 inla 函数进行估计。

```
library(INLA)
wbdata = read.table("wb.txt", header=T)
print(head(wbdata))
formula <- y ~ as.factor(x1) + as.factor(x2) + as.factor(x3)
+ as.factor(x4) + as.factor(x5) + as.factor(x6) + as.factor(x7) +
as.factor(x8)
nm <- inla(formula, data=wbdata, family="nbinomial")
summary(nm)
```

模型的参数估计见表 1 所示。<http://www.r-inla.org> 给出了 R-INLA 的相关说明。

### 4 结论

INLA 方法是一种灵活易用的计算工具。在数据量越来越大的今天,计算开销是重点考虑的问题之一,算法的计算速度显得更加重要,INLA 方法在保证计算精度的同时,具有突出的快速计算能力。由于贝叶斯方法在机器学习领域的重要应用,INLA 方法除了对这一大类隐高斯模型进行贝叶斯推断,在机器学习领域也必然会展现其优势。

表 1 负二项回归基于 INLA 方法的贝叶斯推断

	mean	sd	0.025quant	0.5quant	0.975quant	mode
Intercept	2.9821	0.1131	2.7633	2.981	3.2071	2.9787
as.factor(x <sub>1</sub> )1	-0.487	0.2564	-0.9673	-0.4956	0.0416	-0.5132
as.factor(x <sub>1</sub> )2	-0.1811	0.1085	-0.3953	-0.1807	0.0308	-0.1799
as.factor(x <sub>1</sub> )3	0.0585	0.1388	-0.2118	0.0577	0.3329	0.0562
as.factor(x <sub>1</sub> )4	0.3789	0.1285	0.1275	0.3786	0.6318	0.378
as.factor(x <sub>2</sub> )1	0.464	0.0877	0.2918	0.464	0.636	0.464
as.factor(x <sub>2</sub> )2	2.1568	1.2541	0.0905	2.0105	4.9795	1.679
as.factor(x <sub>3</sub> )1	0.1829	0.1008	-0.0163	0.1833	0.3795	0.1842
as.factor(x <sub>4</sub> )1	0.3565	0.1542	0.0614	0.3538	0.6671	0.3484
as.factor(x <sub>5</sub> )1	-0.1704	0.0898	-0.3455	-0.1709	0.0071	-0.1718
as.factor(x <sub>6</sub> )1	-0.9897	0.18	-1.3336	-0.993	-0.6265	-0.9997
as.factor(x <sub>7</sub> )1	0.3063	0.1396	0.0355	0.3052	0.5834	0.3029
as.factor(x <sub>8</sub> )1	-0.1905	0.0834	-0.3543	-0.1905	-0.0267	-0.1906

#### 参考文献:

- [1]Rue H, Martino S, Chopin N. Approximate Bayesian Inference for latent Gaussian Models Using Integrated Nested Laplace Approximations[J]. Journal of the Royal Statistical Society, Series B, 2009,71(2).
- [2]Martino S, et al. Estimating Stochastic Volatility Models Using Integrated Nested Laplace Approximations[J]. European Journal of Finance, 2011, 17(7).
- [3]Fong Y, Rue H, Wakefield J. Bayesian Inference for Generalized Linear Mixed Models[J]. Biostatistics, 2010, 11(3).
- [4]Yu R Y, Rue H. Bayesian Inference for Additive Mixed Quantile Regression Models[J]. Computational Statistics & Data Analysis, 2011, 55(1).
- [5]Natário I, Oliveira M M, Susete M.Using INLA to Estimate a Highly Dimensional Spatial Model for Forest Fires in Portugal[M]. Berlin: Springer, 2014.
- [6]Gómez-Rubio V, Bivand R, Rue H.A New Latent Class to Spatial Econometrics Models With Integrated Nested Laplace Approximations [J]. Procedia Environmental Sciences, 2015,(27).
- [7]茆诗松,王静龙,濮晓龙.高等数理统计[M].北京:高等教育出版社,2000.
- [8]Paul M, Riebler A, Bachmann L M, et al. Bayesian Bivariate Meta-Analysis of Diagnostic Test Studies Using Integrated Nested Laplace Approximations[J]. Statistics in Medicine, 2010, 29(12).
- [9]Grilli L, Metelli S, Rampichini C. Bayesian Estimation With Integrated Nested Laplace Approximation for Binary Logit Mixed Models[J]. Journal of Statistical Computation and Simulation, 2015,85(13).
- [10]Schrödle B, Held L, Riebler A, et al. Using Integrated Nested Laplace Approximations for the Evaluation Of Veterinary Surveillance Data From Switzerland: A Case-Study[J]. Applied Statistics, 2011, 60(2).
- [11]Taylor B M, Diggle P J. INLA or MCMC? A Tutorial and Comparative Evaluation for Spatial Prediction in Log-Gaussian Cox Processes[J]. Journal of Statistical Computation & Simulation, 2012, 84(10).

(责任编辑/浩 天)