



POLITECNICO
MILANO 1863

MATHEMATICAL ENGINEERING A.Y. 2022-23
BAYESIAN STATISTICS(052499) - PROF. ALESSANDRA GUGLIELMI
TUTORING TEAM: RAFFAELE ARGIENTO, SIRIO LEGRAMANTI, LUCIA PACI

BAYESIAN MIXTURE MODEL FOR ENVIRONMENTAL APPLICATION

TRANSCRIPTION OF THE MODEL
29th January 2023

P. Bogani, P. Botta, S. Caresana, R. Carrara, G. Corbo, L. Mainini

$$\begin{aligned} \mathbf{y}_i &= \mathbf{Z}\boldsymbol{\alpha}_i + \boldsymbol{\theta}_i + \boldsymbol{\epsilon}_i \quad \text{for } i = 1, 2, \dots, n, \\ \boldsymbol{\epsilon}'_i &= (\epsilon_{i1}, \dots, \epsilon_{iT}) \sim N_T(\mathbf{0}, \sigma_{\epsilon_i}^2 \mathbf{I}), \\ \theta_{it} &= \rho_i \theta_{i,t-1} + \nu_{it} \quad \text{with } \nu_{i1} \sim N(0, \sigma_i^2), \\ \gamma_i &= (\sigma_i^2, \rho_i), \\ \boldsymbol{\theta}_i &\sim N_T(\mathbf{0}, \mathbf{R}(\gamma_i)) \quad \text{with } R_{ls}(\gamma_i) = \sigma_i^2 \rho_i^{|l-s|} \text{ and } l, s = 1, \dots, T. \end{aligned}$$

1 Priors distributions

$$\begin{aligned} \boldsymbol{\alpha}_i &\stackrel{\text{iid}}{\sim} N_p(\mathbf{0}, \boldsymbol{\Sigma}_\alpha) \quad \text{with } \boldsymbol{\Sigma}_\alpha = \text{diag}(\sigma_{\alpha_1}, \dots, \sigma_{\alpha_p}), \\ \sigma_{\epsilon_i}^2 &\sim \text{IGa}(c_0^\epsilon, c_1^\epsilon), \quad \sigma_{\alpha_k}^2 \sim \text{IGa}(c_0^\alpha, c_1^\alpha), \\ \gamma_i \mid P &\stackrel{\text{iid}}{\sim} P, \text{ for } i = 1, \dots, n \text{ with } P \sim \text{Dir}(a^P, P_0), \\ p_0(\gamma_i) &= p_0(\sigma_i^2) \times p_0(\rho_i) \\ p_0(\sigma_i^2) &= \text{IGa}(a, b) \quad p_0(\rho_i) = \text{Beta}_{[-1,1]}(c, d), \end{aligned}$$

for $i = 1, \dots, n$, where $\text{Beta}_{[-1,1]}(c, d)$ is a beta distribution with domain in $[-1, 1]$ and hyper-parameters c and d .

2 Posterior distributions

In the following computations, we'll make use of some properties related to multivariate Gaussian distribution that we report here.

Theorem 2.1. *Given a marginal Gaussian distribution for \mathbf{x} and a conditional Gaussian distribution for \mathbf{y} given \mathbf{x} in the form*

$$\begin{aligned} p(\mathbf{x}) &= \mathcal{N}(\mathbf{x} \mid \boldsymbol{\mu}, \boldsymbol{\Lambda}^{-1}), \\ p(\mathbf{y} \mid \mathbf{x}) &= \mathcal{N}(\mathbf{y} \mid \mathbf{A}\mathbf{x} + \mathbf{b}, \mathbf{L}^{-1}), \end{aligned}$$

the marginal distribution of \mathbf{y} and the conditional distribution of \mathbf{x} given \mathbf{y} are given by

$$p(\mathbf{y}) = \mathcal{N}(\mathbf{y} \mid \mathbf{A}\boldsymbol{\mu} + \mathbf{b}, \mathbf{L}^{-1} + \mathbf{A}\boldsymbol{\Lambda}^{-1}\mathbf{A}^T), \quad (1)$$

$$p(\mathbf{x} \mid \mathbf{y}) = \mathcal{N}(\mathbf{x} \mid \boldsymbol{\Sigma} \{ \mathbf{A}^T \mathbf{L}(\mathbf{y} - \mathbf{b}) + \boldsymbol{\Lambda}\boldsymbol{\mu} \}, \boldsymbol{\Sigma}), \quad (2)$$

where

$$\boldsymbol{\Sigma} = (\boldsymbol{\Lambda} + \mathbf{A}^T \mathbf{L} \mathbf{A})^{-1}.$$

A proof of the previous properties can be found in Appendix B (Pag 91-92 of [1]).

If we let $\boldsymbol{\alpha}' = (\boldsymbol{\alpha}'_1, \dots, \boldsymbol{\alpha}'_n)$, $\boldsymbol{\theta}' = (\boldsymbol{\theta}'_1, \dots, \boldsymbol{\theta}'_n)$, and $\boldsymbol{\sigma}'_\epsilon = (\sigma_{\epsilon_1}^2, \dots, \sigma_{\epsilon_n}^2)$ then the likelihood function is given by

$$f(\mathbf{y} \mid \boldsymbol{\alpha}, \boldsymbol{\theta}, \boldsymbol{\sigma}_\epsilon) = \prod_{i=1}^n N_T(\mathbf{y}_i \mid \mathbf{Z}\boldsymbol{\alpha}_i + \boldsymbol{\theta}_i, \sigma_{\epsilon_i}^2 \mathbf{I}).$$

With this notation, we denote that $\mathbf{y}_i \mid \boldsymbol{\alpha}_i, \boldsymbol{\theta}_i, \sigma_{\epsilon_i}^2$ are distributed as an T -variate multivariate normal distribution with mean vector $\mathbf{Z}\boldsymbol{\alpha}_i + \boldsymbol{\theta}_i$ and variance-covariance matrix $\sigma_{\epsilon_i}^2 \mathbf{I}$. We will continue using this notation from now on.

Using Property 1 of Theorem 2.1:

$$f(\mathbf{y}_i \mid \boldsymbol{\alpha}_i, \sigma_{\epsilon_i}^2, \boldsymbol{\Sigma}_\alpha, \gamma_i) = N_p(\mathbf{y}_i \mid \mathbf{Z}\boldsymbol{\alpha}_i, \mathbf{W}_i), \quad (3)$$

with matrices $\mathbf{W}_i = \sigma_{\epsilon_i}^2 \mathbf{I} + \mathbf{R}(\gamma_i)$.

The conditional posterior distribution of $\boldsymbol{\alpha}_i$ becomes

$$f(\boldsymbol{\alpha}_i \mid \mathbf{y}, \sigma_{\epsilon_i}^2, \boldsymbol{\Sigma}_\alpha) = N_p(\boldsymbol{\alpha}_i \mid \boldsymbol{\mu}_\alpha, \mathbf{V}_\alpha), \quad (4)$$

for $i = 1, \dots, n$, where $\boldsymbol{\mu}_\alpha = \mathbf{V}_\alpha \mathbf{Z}' \mathbf{W}_i^{-1} \mathbf{y}_i$ and $\mathbf{V}_\alpha = (\mathbf{Z}' \mathbf{W}_i^{-1} \mathbf{Z} + \boldsymbol{\Sigma}_\alpha^{-1})^{-1}$ with matrices

$$\mathbf{W}_i = \sigma_{\epsilon_i}^2 \mathbf{I} + \mathbf{R}(\gamma_i),$$

of dimensions $T \times T$.

The conditional posterior distribution for the variances $\sigma_{\epsilon_i}^2, i = 1, \dots, n$, and $\sigma_{\alpha_k}^2, k = 1, \dots, p$ given the data and the rest of the parameters are all conditionally conjugate. The conditional posterior distribution for $\sigma_{\epsilon_i}^2$ has the form

$$f(\sigma_{\epsilon_i}^2 \mid \mathbf{y}, \text{rest}) = \text{IGa}\left(\sigma_{\epsilon_i}^2 \mid c_0^\epsilon + \frac{T}{2}, c_1^\epsilon + \frac{1}{2}\mathbf{M}_i'\mathbf{M}_i\right), \quad (5)$$

where $\mathbf{M}_i = \mathbf{y}_i - \mathbf{Z}\boldsymbol{\alpha}_i - \boldsymbol{\theta}_i$, for $i = 1, \dots, n$.

The conditional posterior distribution for $\sigma_{\alpha_j}^2$ has the form

$$f(\sigma_{\alpha_j}^2 \mid \mathbf{y}, \text{rest}) = \text{IGa}\left(\sigma_{\alpha_j}^2 \mid c_0^\alpha + \frac{n}{2}, c_1^\alpha + \frac{1}{2}\sum_{i=1}^n \alpha_{ij}^2\right), \quad (6)$$

for $j = 1, 2, \dots, p$.

The conditional posterior distribution for $\boldsymbol{\theta}_i$ has the form

$$f(\boldsymbol{\theta}_i \mid \mathbf{y}, \sigma_{\epsilon_i}^2, \mathbf{R}(\gamma_i)) = \text{N}_T(\boldsymbol{\theta}_i \mid \boldsymbol{\mu}_\theta, \mathbf{S}_\theta), \quad (7)$$

for $i = 1, \dots, n$, where $\mathbf{S}_\theta = \left((\sigma_{\epsilon_i}^2 \mathbf{I})^{-1} + \mathbf{R}(\gamma_i)^{-1}\right)^{-1}$ and $\boldsymbol{\mu}_\theta = \mathbf{S}_\theta (\sigma_{\epsilon_i}^2 \mathbf{I})^{-1} (\mathbf{y}_i - \mathbf{Z}\boldsymbol{\alpha}_i)$.

Fixing i , we recall that γ_{-i} denotes the set of all γ 's excluding the i^{th} element. $\gamma_{j,i}^*$'s denote the unique values in γ_{-i} , each occurring with frequency $n_{j,i}^*$, with $j = 1, \dots, m_i$, where m_i represents the number of unique values in γ_{-i} .

The posterior distribution for γ is characterized by the formulation of algorithm 8 by Neal [2], using n_{aux} auxiliary variables to prevent the fact that the prior are not conjugate. *da scrivere meglio sta parte*

We can rewrite the model as

$$\begin{aligned} y_i \mid c_i, \boldsymbol{\gamma} &\sim f(\gamma_{c_i}) \\ c_i \mid \mathbf{p} &\sim \text{Discrete}(p_1, \dots, p_K) \\ \gamma_c &\sim P_0 \\ \mathbf{p} &\sim \text{Dirichlet}(a^p/K, \dots, a^p/K). \end{aligned}$$

Here, c_i indicates which "latent class" is associated with observation y_i . The mixing proportions for the classes, $\mathbf{p} = (p_1, \dots, p_K)$, are given a symmetric Dirichlet prior, with concentration parameter written as a^p/K , so that it approaches zero as K goes to infinity. $f(\gamma_{c_i})$ represents the density function:

$$f(\mathbf{y}_i \mid \gamma_{c_i}, \boldsymbol{\alpha}_i, \sigma_{\epsilon_i}^2) = \text{N}_T(\mathbf{y}_i \mid \mathbf{Z}\boldsymbol{\alpha}_i, \sigma_{\epsilon_i}^2 \mathbf{I} + \mathbf{R}(\gamma_{c_i})),$$

due to the theorem 2.1.

We can now update c_i by sampling from its conditional distribution given y_i and the parameters of all existing and empty clusters. Specifically,

For $i = 1, \dots, n$, let m_i the number of distinct c_i for $j \neq i$ and $h = m_i + n_{\text{aux}}$.

If $c_i = c_j$ for some $j \neq i$, draw values independently from P_0 for those γ_c^* for which $m_i < c \leq h$. If $c_i \neq c_j$ for all $j \neq i$, let c_i have the label $m_i + 1$, and draw values independently from P_0 for those γ_c^* for which $m_i + 1 < c \leq h$. Draw a new value for c_i from $\{1, \dots, h\}$ using the following probabilities:

$$\mathbb{P}[c_i = c \mid \mathbf{c}_{-i}, \mathbf{y}_i, \gamma_1^*, \dots, \gamma_h^*] = \begin{cases} b \cdot \frac{n_{c,i}^*}{n - 1 + a^p} \cdot f(y_i \mid \gamma_{c_i}^*, \boldsymbol{\alpha}_i, \sigma_{\epsilon_i}), & \text{for } 1 \leq c \leq m_i \\ b \cdot \frac{a^p / n_{\text{aux}}}{n - 1 + a^p} \cdot f(y_i \mid \gamma_{c_i}^*, \boldsymbol{\alpha}_i, \sigma_{\epsilon_i}), & \text{for } m_i < c \leq h \end{cases}$$

where $n_{c,i}^*$ is the number of c_j for $j \neq i$ that are equal to c , and b is the appropriate normalizing constant. The values of $\gamma_{c_i}^*$ where $m_i < c \leq h$ are drawn independently from P_0 .

References

- [1] Christopher M Bishop and Nasser M Nasrabadi. *Pattern recognition and machine learning*. Vol. 4. 4. Springer, 2006.
- [2] Radford M Neal. ‘Markov chain sampling methods for Dirichlet process mixture models’. In: *Journal of computational and graphical statistics* 9.2 (2000), pp. 249–265.