

Modèles à variables latentes pour l'écologie et la biologie

Examen

Note that a few useful formula are given at the end of the document.

We are interested in modeling the count of individuals of the same species on several sites. Let i be the geographical site ($i = 1, \dots, N$). Each site i is visited n_i times. Y_{ij} is the number of individuals observed at site i observation j .

Part 1. Frequentist estimation

We consider the following model:

$$\begin{cases} Y_{ij}|Z_i & \sim_{ind} \mathcal{P}(e^{Z_i}) \\ Z_i & \sim_{i.i.d} \mathcal{N}(\mu, \sigma^2). \end{cases} \quad (1)$$

(Z_i) introduces a variability due to the heterogeneity of the sites.

In what follows, we will use the following notations:

$$\begin{aligned} \mathbf{Y} &:= (Y_{ij})_{i=1,\dots,N,j=1,\dots,n_i} \\ \mathbf{Z} &:= (Z_i)_{i=1,\dots,N} \\ \theta &:= (\mu, \sigma^2) \end{aligned}$$

1. Give the expression of the complete log likelihood $\log p(\mathbf{Y}, \mathbf{Z}; \theta)$.
2. Are you able to give a close form expression of the likelihood of the observations \mathbf{Y} and of the conditional distribution $p(\mathbf{Z}|\mathbf{Y}; \theta)$? (Explain why)
3. Recall the general principle of the Variational EM algorithm (lower bound, VE and M step, algorithm...) [[Question de cours](#)]

We propose to approximate the conditional distributions $p(Z_i|\mathbf{Y})$ in the Gaussian family:

$$\tilde{q}(\mathbf{z}) = \prod_{i=1}^N \tilde{q}_i(z_i) \quad \text{where} \quad \tilde{q}_i(z_i) = f_{\mathcal{N}(\tilde{\mu}_i, \tilde{\omega}_i^2)}(z_i) \quad (2)$$

4. Prove that the entropy of \tilde{q} is

$$\mathcal{H}(\tilde{q}) = \sum_{i=1}^N \left(\frac{1}{2} \log(2\pi) + \frac{1}{2} \log \tilde{\omega}_i^2 + 1 \right)$$

5. Derive the expression of the lower bound as a function of θ and $\tau = (\tilde{\mu}_i, \tilde{\omega}_i^2)_{i=1,\dots,N}$

6. Give the expression of $\hat{\theta}$ solution of the M-step
7. Give the equation verified by $\hat{\tau}$ solution of the VE-step. In practice, how do you propose to solve it?
8. Propose an (efficient) initialisation of your VEM algorithm.

Part 2. Bayesian estimation

Assume that each site is described by a collection of p environmental covariates. Let $x_i \in \mathbb{R}^p$ the covariates for site $i = 1, \dots, N$.

We set the following model

$$\begin{cases} Y_{ij}|Z_i & \sim_{ind} \mathcal{P}(e^{Z_i}) \\ Z_i & \sim_{ind} \mathcal{N}(x_i^T \beta, \sigma^2). \end{cases} \quad (3)$$

where $\beta \in \mathbb{R}^p$. In this model, a part of the variability between sites is explained by the covariates.

We assume that the covariate x_1, \dots, x_N are such that the $n \times p$ matrix $X = \begin{pmatrix} x_1 \\ \vdots \\ x_N \end{pmatrix}$ is of full rank. For the sake of simplicity, σ^2 is assumed to be known.

We consider a Bayesian inference for β et set the following prior distribution on β :

$$\beta \sim \mathcal{N}(0_p, \omega^2 \mathbf{I}_p) \quad (4)$$

where 0_p is the null vector of size p , \mathbf{I}_p is the identity matrix of size p . $\omega^2 \in \mathbb{R}^{+*}$.

9. What is the role of ω^2 in the prior distribution (4)?
10. Give the expression of $p(\mathbf{Y}, \mathbf{Z}, \beta; \omega^2, \sigma^2)$ for Model (3) and prior distribution on β (4).

The aim is to propose a method to obtain a sample from $p(\beta, \mathbf{Z}|\mathbf{Y})$. Let h be the iteration number. We propose to sample iteratively:

- a. $\beta^{(h)} \sim p(\beta|\mathbf{Z}^{(h-1)}, \mathbf{Y}; \omega^2, \sigma^2)$
- b. $\mathbf{Z}^{(h)} \sim p(\mathbf{Z}|\beta^{(h)}, \mathbf{Y}; \omega^2, \sigma^2)$

11. [\[Question de cours\]](#) What is the name of the algorithm? Give quickly its properties.
12. Show that the simulation at step [a.] can be performed exactly (give the distribution of β given the latent variable \mathbf{Z} and the data \mathbf{Y}).
13. Are you able to simulate explicitly $\mathbf{Z}^{(h)} \sim p(\mathbf{Z}|\beta^{(h)}, \mathbf{Y}; \omega^2, \sigma^2)$?
14. [\[Question de cours\]](#) Recall the principle of the Metropolis-Hastings algorithm.

15. How can you apply it here at step [b.] of the algorithm?
16. How can you tune the algorithm to reach the adequate convergence rate?

Useful formulae

- $\mathbb{E}[Z^2] = \mathbb{V}(Z) + (\mathbb{E}[Z])^2$
- $\mathcal{H}(q) = \mathbb{E}_{Z \sim q}[\log q(Z)]$
- Assume that $Z \sim \mathcal{N}(\mu, \sigma^2)$ then
 - ★ the density is $f(z) = \frac{1}{\sqrt{2\pi\sigma^2}} e^{-\frac{1}{2\sigma^2}(z-\mu)^2}$
 - ★ $\mathbb{E}(e^Z) = e^{\frac{\sigma^2}{2} + \mu}$
- **Gaussian vector** Let X be a full rank matrix of size $n \times p$ ($n < p$). If

$$\begin{aligned} \mathbf{Z} = (Z_1, \dots, Z_n) &\sim \mathcal{N}_n(X\beta, \sigma^2 \mathbf{I}_n) \\ \text{and } \beta &\sim \mathcal{N}_p(0_p, \omega^2 \mathbf{I}_p) \end{aligned}$$

then

$$\begin{aligned} \beta | \mathbf{Z} &\sim \mathcal{N}_p(\mu^{post}, \Omega^{post}) \\ \text{with } \Omega^{post} &= \left(X^T X + \frac{\sigma^2}{\omega^2} \mathbf{I}_p \right)^{-1} \\ \text{and } \mu^{post} &= \Omega^{post} X^T \mathbf{Z} \end{aligned}$$

where X^T is the transposed matrix of X .