

Eyes : Normal Mixture Model

Le ZHANG & Shijie XU

13 mars 2025

1 Introduction

Cette étude analyse les données fournies par Bowmaker et al. (1985), qui concernent les longueurs d'onde de sensibilité maximale des photorécepteurs d'un singe. L'ensemble de données contient 48 observations, et l'objectif est de modéliser et d'estimer les moyennes des deux groupes, leur proportion et leur variance à l'aide d'un modèle de mélange normal (Normal Mixture Model).

2 Modèle et Méthodologie

Nous supposons que chaque observation y_i provient de l'une des deux distributions normales :

$$y_i \sim \mathcal{N}(\lambda_{T_i}, \tau)$$

où :

- $T_i \sim \text{Categorical}(P)$ représente le groupe auquel appartient chaque observation ($T_i = 1$ ou 2).
- La moyenne du premier groupe est λ_1 , et celle du deuxième groupe est définie comme $\lambda_2 = \lambda_1 + \theta$, $\theta > 0$.
- Les deux groupes partagent une même variance contrôlée par le paramètre de précision τ .

λ_1 , θ , τ , P sont dotés de prioris indépendants dits “non-informatifs”, y compris un priori uniforme pour P sur $(0, 1)$. Ainsi, nous adoptons les lois a priori suivantes :

$$\theta \sim \mathcal{N}(0, \sigma_\theta^2) \cdot \mathbb{I}_{[0, +\infty]}$$

$$\tau \sim \text{Gamma}(\alpha, \beta)$$

$$\lambda_1 \sim \mathcal{N}(0, \sigma_1^2)$$

$$P \sim \text{Dirichlet}(1, 1)$$

Voici les détails mathématiques de l'échantillonnage MCMC.

En introduisant les variables latentes T_i , la vraisemblance complète s'écrit :

$$\pi(\mathbf{y}, \mathbf{T} \mid \lambda_1, \theta, \tau, P) = \prod_{i=1}^N P_{T_i} \cdot \phi_\tau(y_i - \lambda_{T_i})$$

avec $\phi_\tau(y - \lambda)$ la densité d'une loi normale de précision τ :

$$\phi_\tau(y - \lambda) = \sqrt{\frac{\tau}{2\pi}} \exp\left(-\frac{\tau}{2}(y - \lambda)^2\right)$$

L'expression de la loi a posteriori jointe est :

$$\pi(\lambda_1, \theta, \tau, P, \mathbf{T} \mid \mathbf{y}) \propto \pi(\mathbf{y}, \mathbf{T} \mid \lambda_1, \theta, \tau, P) \cdot \pi(\lambda_1) \cdot \pi(\theta) \cdot \pi(\tau) \cdot \pi(P)$$

1. Mettre à jour T_i via une loi de Bernoulli :

$$\mathbb{P}(T_i = 1 \mid y_i, \lambda_1, \theta, \tau, P) = \frac{P_1 \cdot \phi_\tau(y_i - \lambda_1)}{P_1 \cdot \phi_\tau(y_i - \lambda_1) + P_2 \cdot \phi_\tau(y_i - (\lambda_1 + \theta))}$$

On tire alors T_i selon une loi de Bernoulli avec ce paramètre.

2. Mettre à jour P via une loi Beta : Une fois tous les T_i mis à jour, on compte le nombre d'observations par groupe :

$$n_1 = \sum_{i=1}^N 1(T_i = 1), \quad n_2 = N - n_1$$

Puis, la loi conditionnelle de P est :

$$P \sim \text{Dirichlet}(1 + n_1, 1 + n_2)$$

Soit :

$$P_1 \sim \text{Beta}(1 + n_1, 1 + n_2), \quad P_2 = 1 - P_1.$$

3. Mettre à jour (λ_1, θ) via Metropolis-Hastings : Puisque $\lambda_2 = \lambda_1 + \theta$ et que $\theta > 0$, il n'existe pas de mise à jour conjuguée. On utilise donc un échantillonnage Metropolis-Hastings sur (λ_1, θ) :

$$\pi(\lambda_1, \theta \mid T, y_i, \tau) \propto \pi(y_i \mid \lambda_1, \theta, \tau, T) \cdot \pi(\lambda_1) \cdot \pi(\theta)$$

$$\log \pi(\lambda_1, \theta \mid T, y_i, \tau) \propto \underbrace{\log \pi(\lambda_1) + \log \pi(\theta)}_{\text{petit terme}} + \log \pi(y_i \mid \lambda_1, \theta, \tau, T)$$

$$\propto \sum_{i:T_i=1} \left[-\frac{\tau}{2}(y_i - \lambda_1)^2 \right] + \sum_{i:T_i=2} \left[-\frac{\tau}{2}(y_i - (\lambda_1 + \theta))^2 \right]$$

On échantillonne (λ_1, θ) conjointement via Metropolis-Hastings, en imposant $\theta > 0$.

4. Mettre à jour τ via une loi Gamma :

$$\begin{aligned}\pi(\tau \mid \lambda, \theta, T, y) &\propto \pi(y \mid \tau, \lambda, \theta, T) \cdot \pi(\tau) \\ &\propto \prod_{i=1}^N \left[\tau^{\frac{1}{2}} \exp \left(-\frac{\tau}{2} (y_i - \lambda_{T_i})^2 \right) \right] \cdot \tau^{\alpha-1} e^{-\beta\tau} \\ &\propto \tau^{\alpha + \frac{N}{2} - 1} \exp \left[-\tau \left(\beta + \frac{1}{2} \sum_{i=1}^N (y_i - \lambda_{T_i})^2 \right) \right] \\ \tau &\sim \text{Gamma} \left(\alpha + \frac{N}{2}, \beta + \frac{1}{2} \sum_{i=1}^N (y_i - \lambda_{T_i})^2 \right)\end{aligned}$$

5. Mettre à jour σ via $\sigma = 1/\sqrt{\tau}$.

3 Résultats et Analyse

Après 10 000 itérations avec 1 000 de burn-in, nous obtenons les estimations postérieures suivantes :

| Parameter | Mean | SD | MC_error | 2.5% | Median | 97.5% | start | Sample |
|-----------|----------|--------|----------|----------|----------|----------|-------|--------|
| P[1] | 0.6018 | 0.0826 | 0.000870 | 0.4392 | 0.6036 | 0.7586 | 1001 | 10000 |
| P[2] | 0.3982 | 0.0826 | 0.000870 | 0.2414 | 0.3964 | 0.5608 | 1001 | 10000 |
| lambda[1] | 536.7021 | 0.9058 | 0.009548 | 534.8288 | 536.6955 | 538.4622 | 1001 | 10000 |
| lambda[2] | 548.9633 | 1.1541 | 0.012165 | 546.6067 | 548.9466 | 551.3026 | 1001 | 10000 |
| sigma | 3.7330 | 0.5412 | 0.005704 | 2.9152 | 3.6545 | 5.0253 | 1001 | 10000 |

FIGURE 1 – Estimations postérieures après 10 000 itérations

Ces résultats sont cohérents avec le document fourni. Les incertitudes sont faibles, indiquant une convergence satisfaisante du MCMC.

Ce modèle de mélange normal permet donc d'identifier deux sous-groupes distincts dans les données, et les estimations obtenues sont robustes et fiables. 请再写一点，我有点累。

Annexe

A1 DAG

Voici la DAG correspondant :

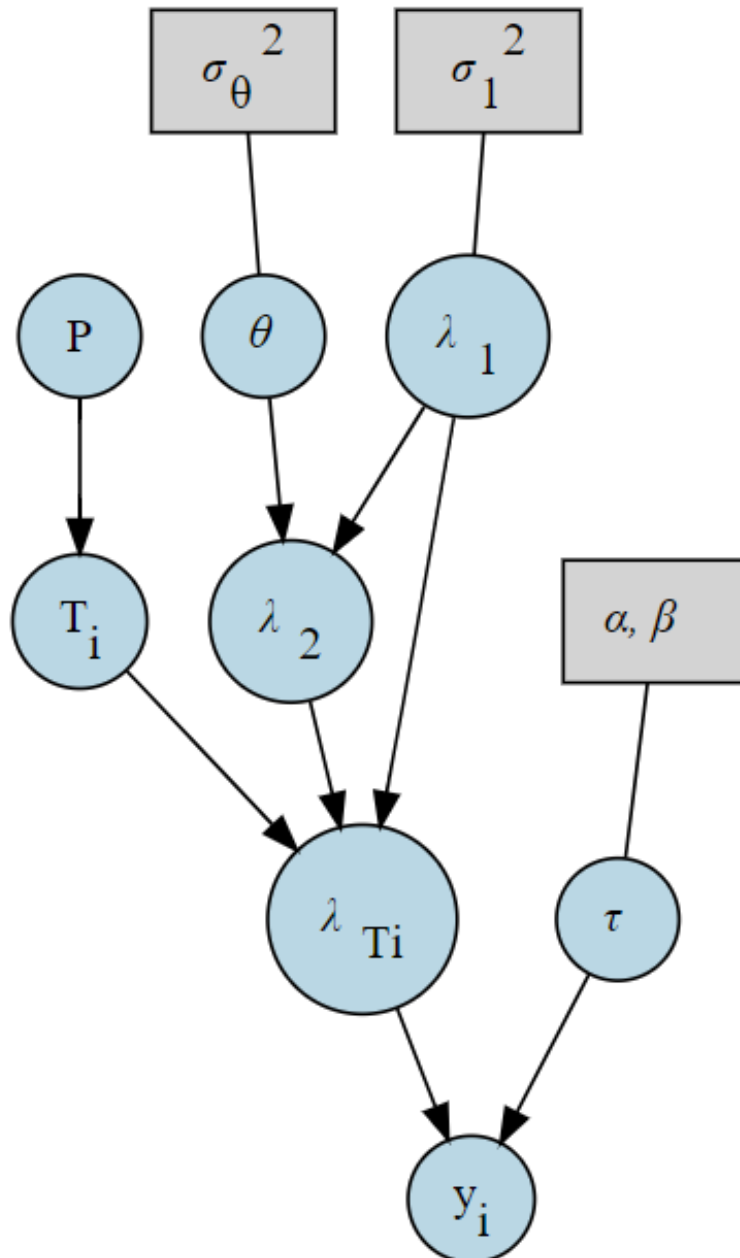


FIGURE 2 – DAG correspondante

A2 code

```
1  ```{r}
2  N <- 48
3  y <- c(529, 530, 532, 533.1, 533.4, 533.6, 533.7, 534.1, 534.8, 535.3,
4        535.4, 535.9, 536.1, 536.3, 536.4, 536.6, 537, 537.4, 537.5, 538.3,
5        538.5, 538.6, 539.4, 539.6, 540.4, 540.8, 542, 542.8, 543, 543.5,
6        543.8, 543.9, 545.3, 546.2, 548.8, 548.7, 548.9, 549, 549.4, 549.9,
7        550.6, 551.2, 551.4, 551.5, 551.6, 552.8, 552.9, 553.2)
8
9  alpha <- c(1,1)    # Dirichlet(1,1)
10
11  ##MCMC
12  n_iter <- 10000
13  burn_in <- 1000
14  thin <- 1
15  set.seed(123)
16
17  ## initiation
18  lambda1 <- 535
19  theta <- 5
20  tau <- 1/10    # sigmasq = 10 donc tau = 0.1
21  p <- c(0.5, 0.5)
22
23  T_ <- sample(1:2, N, replace=TRUE)
24
25  ##MCMC sampling
26  samples_lambda1 <- numeric(n_iter)
27  samples_theta <- numeric(n_iter)
28  samples_tau <- numeric(n_iter)
29  samples_p1 <- numeric(n_iter)
30
31
32  ```
33
34  ```{r}
35  sigma_lambda <- 0.5
36  sigma_theta <- 0.2
37
38  accept_count <- 0
39  adjust_interval <- 100
40
41  for(iter in 1:n_iter) {
42
43    ## update T[i]
44    lam2 <- lambda1 + theta
```

```

45 for(i in 1:N) {
46   prob1 <- p[1] * dnorm(y[i], lambda1, sqrt(1/tau))
47   prob2 <- p[2] * dnorm(y[i], lam2, sqrt(1/tau))
48   T_[i] <- sample(1:2, 1, prob=c(prob1, prob2))
49 }
50
51 ## update P
52 n1 <- sum(T_ == 1)
53 p1_new <- rbeta(1, 1 + n1, 1 + (N - n1))
54 p <- c(p1_new, 1 - p1_new)
55
56 ## update Metropolis (lambda1, theta)
57 lam1_prop <- lambda1 + rnorm(1, mean=0, sd=sigma_lambda)
58 theta_prop <- abs(theta + rnorm(1, mean=0, sd=sigma_theta))
59
60 log_accept_ratio <- sum(dnorm(y[T_==1], mean=lam1_prop, sd=sqrt(1/tau),
61   log=TRUE)) +
62   sum(dnorm(y[T_==2], mean=lam1_prop + theta_prop, sd=sqrt(1/tau), log=
63     TRUE)) -
64   sum(dnorm(y[T_==1], mean=lambda1, sd=sqrt(1/tau), log=TRUE)) -
65   sum(dnorm(y[T_==2], mean=lambda1 + theta, sd=sqrt(1/tau), log=TRUE))
66
67 if(log(runif(1)) < log_accept_ratio) {
68   lambda1 <- lam1_prop
69   theta <- theta_prop
70   accept_count <- accept_count + 1
71 }
72
73 ## (d) update tau
74 ssq <- sum((y[T_==1] - lambda1)^2) + sum((y[T_==2] - (lambda1 + theta))^2)
75
76 tau <- rgamma(1, shape=0.001 + N/2, rate=0.001 + 0.5 * ssq)
77
78 ## save the samples
79 samples_lambda1[iter] <- lambda1
80 samples_theta[iter] <- theta
81 samples_tau[iter] <- tau
82 samples_p1[iter] <- p[1]
83
84 if(iter %% adjust_interval == 0) {
85   accept_rate <- accept_count / adjust_interval # calculate accept rate
86   if(accept_rate < 0.2) {
87     sigma_lambda <- sigma_lambda * 0.9
88     sigma_theta <- sigma_theta * 0.9
89   } else if(accept_rate > 0.5) {
90     sigma_lambda <- sigma_lambda * 1.1

```

```

88     sigma_theta <- sigma_theta * 1.1
89   }
90   accept_count <- 0
91 }
92 }
93
94 ## analyse the result
95 burned_samples <- (burn_in+1):n_iter
96 cat("Posterior_Mean_Estimates:\n")
97 cat("lambda1=", mean(samples_lambda1[burned_samples]), "\n")
98 cat("theta_111=", mean(samples_theta[burned_samples]), "\n")
99 cat("tau_11111=", mean(samples_tau[burned_samples]), "\n")
100 cat("p1_111111=", mean(samples_p1[burned_samples]), "\n")
101 ```
102
103 ```{r}
104
105 posterior_summary <- function(samples, name) {
106   mean_val <- mean(samples)
107   sd_val <- sd(samples)
108   mc_error <- sd_val / sqrt(length(samples))
109   quantiles <- quantile(samples, probs = c(0.025, 0.5, 0.975))
110
111   cat(sprintf("%-10s_8.4f_8.4f_10.6f_8.4f_8.4f_8.4f_6d_6d\n",
112               name, mean_val, sd_val, mc_error,
113               quantiles[1], quantiles[2], quantiles[3], burn_in+1, n_iter))
114 }
115
116 cat(sprintf("%-10s_8s_8s_10s_8s_8s_8s_6s_6s\n",
117             "Parameter", "Mean", "SD", "MC_error", "2.5%", "Median", "97.5%",
118             "Start", "Sample"))
119
120 cat(rep("-", 70), "\n")
121
122 posterior_summary(samples_p1[burned_samples], "P[1]")
123 posterior_summary(1 - samples_p1[burned_samples], "P[2]")
124 posterior_summary(samples_lambda1[burned_samples], "lambda[1]")
125 posterior_summary((samples_lambda1 + samples_theta)[burned_samples], "
    lambda[2]")
126 posterior_summary((1 / sqrt(samples_tau))[burned_samples], "sigma") #
    sigma = 1/sqrt(tau)
127 ```

```