



Université de Montpellier

Projet M1 SSD

Un modèle pour les nids de mouettes

Rédigé par

CARVAILLO Thomas

CÔME Olivier

PRALON Nicolas

Encadrante : Elodie BRUNEL-PICCININI

1^{er} mars 2022

Table des matières

| | |
|--|----------|
| Introduction | 2 |
| 1 Un peu de théorie | 3 |
| 1.1 Modélisation du problème | 3 |
| 1.2 Un cas élémentaire | 4 |
| 1.3 Le cas réel | 6 |
| Bibliographie | 7 |
| A Annexe | 8 |

Introduction

Nos amis les ornithologues aiment bien se balader en slibard et compter les nids de mouettes, il faut bien s'occuper pendant le chômage. Hélas, ce sont des flemmards, ils mesurent la taille des nids mais n'attendent pas que le piaf reviennent pour savoir à quel espèce il appartient. Heureusement, nous sommes là pour combler l'incompétence de ces gueux, grâce à nos techniques du turfu nous allons estimer BLABLABLA.

Chapitre 1

Un peu de théorie

1.1 Modélisation du problème

Nous allons pour commencer donner une première définition, qui est au coeur du présent projet.

Définition 1 (Loi de mélange). On appelle loi de mélange toute loi dont la densité s'écrit sous la forme d'une combinaison linéaire de diverses densités. C'est-à-dire que si l'on se donne n variables aléatoires X_1, \dots, X_n de densité respective $f_1(x), \dots, f_n(x)$, alors est appelée loi de mélange toute variable aléatoire X dont la densité f s'exprime sous la forme

$$f(x) := \sum_{i=1}^n \alpha(i) f_i(x), \alpha(i) \in \mathbb{R}$$

Continuons, afin de modéliser commodément le problème, par introduire les variables aléatoires suivantes :

- ✂ La variable aléatoire X , modélisant la taille des nids
- ✂ La variable aléatoire Y , décrivant la taille du nid d'une espèce donnée
- ✂ Et Z , la variable aléatoire représentant l'espèce de mouette qui a construit le nid

Hypothèse 1. Nous supposons que la taille des nids d'une espèce j (*i.e.* X conditionnellement à $(Z = j)$) suit une loi normale $\mathcal{N}(\mu_j, v_j)$.

Proposition 1. La variable Z est discrète et à valeur dans un sous-ensemble fini de \mathbb{N} , elle suit donc une loi

$$\sum_{j=1}^J \alpha(j) \delta_j$$

où J représente le nombre d'espèce de mouettes considéré et les $\alpha(j)$ sont des réels, positifs stricts, représentant la proportion de nids de l'espèce j , tels que $\sum_{j=1}^J \alpha(j) = 1$.

Il s'ensuit la proposition suivante, qui sera la racine du présent projet.

Proposition 2. La distribution de la taille des nids de mouettes, *i.e.* X , admet pour densité, au point x et par rapport à la mesure de Lebesgue sur \mathbb{R} , la fonction f_θ définie comme suit

$$f_\theta(x) = \sum_{j=1}^J \alpha(j) \gamma_{\mu_j, v_j}(x)$$

Démonstration. On vérifie que l'on obtient bien une densité de probabilité, la forme de cette dernière étant la conséquence directe de la définition de la variable aléatoire X :

$$\int_{\mathbb{R}} f_\theta(x) dx = \int_{\mathbb{R}} \sum_{j=1}^J \alpha(j) \gamma_{\mu_j, v_j}(x) dx = \sum_{j=1}^J \alpha(j) \int_{\mathbb{R}} \gamma_{\mu_j, v_j}(x) dx = \sum_{j=1}^J \alpha(j) = 1$$



Remarque 1. Il est intéressant de noter que si f_θ est bien une densité, cela vient du fait que la combinaison à l'origine du mélange de loi est une combinaison *convexe*. ET SINON ?

Exemple 1. ICI une simulation ?

Le but de ce projet sera d'étudier des méthodes permettant l'estimation des divers paramètres de cette densité. Nous dénoterons par $\theta := (\alpha_j, \mu_j, v_j)_{1 \leq j \leq J}$ les vecteurs des ces dits paramètres. Pour cela, nous nous placerons sous l'hypothèse suivante

Hypothèse 2. Soit $\Theta := \{\theta = (\alpha_j, \mu_j, v_j)_{1 \leq j \leq J} \text{ tels que } \alpha_j > 0 \forall j \in \llbracket 1, J \rrbracket \text{ et } \sum_{j=1}^J \alpha_j = 1\}$. Soient X_1, \dots, X_n un échantillon iid de même loi que X . On supposera qu'il existe un $\theta \in \Theta$ tel que les données récoltées, ici les tailles des nids, soient la réalisation du précédent échantillon.

Notation 1 (Densités). Le vecteur θ ayant été dûment introduit, nous noterons

1. $g_\theta(z) = \alpha(z)$ la densité de la variable aléatoire Z
2. $f_\theta(x|Z = j) := \gamma_{\mu_j, v_j}(x)$ la densité de la loi de X sachant Z , i.e. de la variable aléatoire Y

qui sont, respectivement, contre la mesure de comptage sur \mathbb{N} , et par rapport à la mesure de Lebesgue sur \mathbb{R} ,

Introduisons deux dernières densités, qui nous seront fort utile quant à l'expression des Log-vraisemblances conditionnelles :

Proposition 3 (Densité de Z sachant X). La densité de la loi de Z sachant X est donnée par

$$g_\theta(z|X = x) = \frac{\alpha(z)\gamma_{\mu_z, v_z}(x)}{\sum_{j=1}^J \alpha(j)\gamma_{\mu_j, v_j}(x)}, \text{ pour } z \in \{1, \dots, J\}$$

Démonstration. En effet, BLABLABLA Nicolas le fera



Proposition 4 (Densité du vecteur (X, Z)). Soit le vecteur aléatoire (X, Z) ; sa densité nous est donnée par

$$h_\theta(x, z) := \alpha(z)\gamma_{\mu_z, v_z}(x)$$

Démonstration. En effet, si on dénote par F la fonction de répartition de ce vecteur aléatoire, on obtient

$$\begin{aligned} F(x, z) &= \mathbb{P}(X \leq x, Z = z) \\ &= \mathbb{P}(X \leq x | Z = z) \times \mathbb{P}(Z = z) \\ &= \int_{-\infty}^x \gamma_{\mu_z, v_z}(t) dt \times \alpha(z) \\ &= \int_{-\infty}^x \int_{\mathbb{N}} \alpha(z)\gamma_{\mu_z, v_z}(t) dt d\delta_z \end{aligned}$$

La densité de (X, Z) s'ensuit.



Nous allons dès à présent nous intéresser à l'estimation de ces paramètres.

1.2 Un cas élémentaire

Regardons dans un premier temps un cas simplifié, un cas ne décrivant pas la réalité des observations mais qui a le mérite de constituer une agréable entrée en matière.

Nous supposons ici qu'ont été relevés simultanément et les mesures des tailles des nids et l'espèce de mouette qui l'a construit. Le modèle ici considéré est donc composé des couple (X_i, Z_i) , $i \in \llbracket 1, n \rrbracket$. On considérera dès lors la fonction de densité $h_\theta(x, z)$.

L'estimation des divers paramètres est alors élémentaire, en témoigne les propositions suivante :

Proposition 5 (Fonction de Log-vraisemblance). La Log-vraisemblance du modèle s'écrit

$$\mathcal{L}_\theta(X_1, \dots, X_n, Z_1, \dots, Z_n) = \sum_{j=1}^J \#A_j \ln(\alpha(j)) + \sum_{j=1}^J \sum_{i \in A_j} \ln(\gamma_{\mu_j, v_j}(X_i))$$

où les A_j sont définis par $A_j := \{i \in \llbracket 1, n \rrbracket \text{ tels que } Z_i = j\}$ $\bigcup_{j=1}^J A_j = \llbracket 1, n \rrbracket$

Démonstration. La Log-vraisemblance du modèle s'écrit :

$$\begin{aligned} \mathcal{L}_\theta(X_1, \dots, X_n, z_1, \dots, z_n) &= \ln \left(\prod_{i=1}^n \alpha(z_i) \gamma_{\mu_j, v_j}(x_i) \right) \\ &= \sum_{i=1}^n \ln(\alpha(z_i)) + \ln(\gamma_{\mu_j, v_j}(x_i)) \end{aligned}$$

z_i est à valeur dans $\llbracket 1, J \rrbracket$, on partitionne donc $I := \llbracket 1, n \rrbracket$ comme $I = \bigcup_{j=1}^J A_j$ pour obtenir

$$\begin{aligned} \mathcal{L}_\theta(X_1, \dots, X_n, Z_1, \dots, Z_n) &= \sum_{j=1}^J \sum_{i \in A_j} \ln(\alpha(z_i)) + \sum_{j=1}^J \sum_{i \in A_j} \ln(\gamma_{\mu_j, v_j}(x_i)) \\ &= \sum_{j=1}^J \sum_{i \in A_j} \ln(\alpha(j)) + \sum_{j=1}^J \sum_{i \in A_j} \ln(\gamma_{\mu_j, v_j}(x_i)) \\ &= \sum_{j=1}^J \#A_j \ln(\alpha(j)) + \sum_{j=1}^J \sum_{i \in A_j} \ln(\gamma_{\mu_j, v_j}(x_i)) \end{aligned}$$

2

Nous pouvons dès lors maximiser la log-vraisemblance afin d'obtenir les estimateurs souhaités :

Proposition 6 (Estimateurs). *Les estimateurs du maximum de vraisemblance $\hat{\alpha}(j)$ (resp. $\hat{\mu}_j$, et \hat{v}_j) de $\alpha(j)$ (resp. μ_j et v_j) sont donnés par*

$$\begin{aligned} \hat{\alpha}(j) &= \frac{\#A_j}{n} \\ \hat{\mu}_j &= \frac{\sum_{i \in A_j} X_i}{\#A_j} \\ \hat{v}_j &= \frac{\sum_{i \in A_j} (X_i - \hat{\mu}_j)^2}{\#A_j} \end{aligned}$$

Démonstration. Soit $\theta = (\alpha(j), \mu_j, v_j)_{j \in \llbracket 1, J \rrbracket}$. Il s'agit de déterminer

$$\operatorname{argmax}_{\theta \in \mathbb{R}^{3n}, \sum_{j=1}^J \alpha(j)=1} \left(\sum_{j=1}^J \#A_j \ln(\alpha(j)) + \sum_{j=1}^J \sum_{i \in A_j} \ln(\gamma_{\mu_j, v_j}(x_i)) \right)$$

Nous avons donc à résoudre un programme de minimisation d'une fonction convexe sur un convexe avec une contrainte égalité, il est ainsi naturel de faire appel au Lagrangien.

Ce dernier s'écrit

$$\begin{aligned}
L(\theta) &= \sum_{j=1}^J \#A_j \ln(\alpha(j)) + \sum_{j=1}^J \sum_{i \in A_j} \ln(\gamma_{\mu_j, v_j}(x_i)) - \lambda \times \left(\sum_{j=1}^J \alpha(j) - 1 \right) \\
&= \sum_{j=1}^J \#A_j \ln(\alpha(j)) + \sum_{j=1}^J \sum_{i \in A_j} \ln \left(\frac{1}{\sqrt{2\pi v_j}} \exp \left(-\frac{(x_i - \mu_j)^2}{2v_j} \right) \right) - \lambda \times \left(\sum_{j=1}^J \alpha(j) - 1 \right) \\
&= \sum_{j=1}^J \#A_j \ln(\alpha(j)) + \sum_{j=1}^J \sum_{i \in A_j} \left((-1/2) \ln(2\pi v_j) - \frac{(x_i - \mu_j)^2}{2v_j} \right) - \lambda \times \left(\sum_{j=1}^J \alpha(j) - 1 \right)
\end{aligned}$$

Il reste maintenant à résoudre le système suivant, afin d'obtenir le vecteur $\hat{\theta} := (\hat{\alpha}(j), \hat{\mu}_j, \hat{v}_j)_{j \in \llbracket 1, J \rrbracket}$ solution du programme.

$$\begin{cases} \frac{\#A_j}{\hat{\alpha}(j)} - \lambda &= 0 \quad \forall j \in \llbracket 1, J \rrbracket \\ \sum_{i \in A_j} (x_i - \hat{\mu}_j) / \hat{v}_j &= 0 \quad \forall j \in \llbracket 1, J \rrbracket \\ \sum_{i \in A_j} \frac{-0.5 * 2 * \pi}{2\pi \hat{v}_j} + \frac{(x_i - \hat{\mu}_j)^2}{2\hat{v}_j^2} &= 0 \quad \forall j \in \llbracket 1, J \rrbracket \\ \sum_{j=1}^J \hat{\alpha}(j) &= 1 \end{cases}$$

Ceci équivaut à

$$\begin{cases} \frac{\#A_j}{\hat{\alpha}(j)} &= \lambda \quad \forall j \in \llbracket 1, J \rrbracket \\ \sum_{i \in A_j} x_i &= \sum_{i \in A_j} \hat{\mu}_j \quad \forall j \in \llbracket 1, J \rrbracket \\ \sum_{i \in A_j} (x_i - \hat{\mu}_j)^2 &= \sum_{i \in A_j} \hat{v}_j \quad \forall j \in \llbracket 1, J \rrbracket \\ \sum_{j=1}^J \hat{\alpha}(j) &= 1 \end{cases} \Leftrightarrow \begin{cases} \frac{\#A_j}{\hat{\alpha}(j)} &= \lambda \quad \forall j \in \llbracket 1, J \rrbracket \\ \sum_{i \in A_j} \frac{x_i}{\#A_j} &= \hat{\mu}_j \quad \forall j \in \llbracket 1, J \rrbracket \\ \sum_{i \in A_j} \frac{(x_i - \hat{\mu}_j)^2}{\#A_j} &= \hat{v}_j \quad \forall j \in \llbracket 1, J \rrbracket \\ \sum_{j=1}^J \hat{\alpha}(j) &= 1 \end{cases}$$

En sommant les J premières lignes du système, on obtient $\sum_{j=1}^J \#A_j = \sum_{j=1}^J \hat{\alpha}(j) \lambda$, i.e. $\lambda = n$. En injectant ceci dans le précédent système, on obtient finalement ce qui était annoncé :

$$\begin{cases} \hat{\alpha}(j) &= \frac{\#A_j}{n} \quad \forall j \in \llbracket 1, J \rrbracket \\ \hat{\mu}_j &= \sum_{i \in A_j} \frac{x_i}{\#A_j} \quad \forall j \in \llbracket 1, J \rrbracket \\ \hat{v}_j &= \sum_{i \in A_j} \frac{(x_i - \hat{\mu}_j)^2}{\#A_j} \quad \forall j \in \llbracket 1, J \rrbracket \end{cases}$$

1.3 Le cas réel



Bibliographie

Documents qui m'ont l'air pas mal, à voir s'ils seront utiles

<https://www.lpsm.paris/pageperso/rebafka/BookGraphes/algorithmme-em.html>

<https://members.loria.fr/moberger/Enseignement/AVR/Exposes/algo-em.pdf>

Annexe A

Annexe