

TD/TP5: Modèle de capture-recapture

Le contexte typique dans lequel le modèle de capture-recapture est utilisé est celui d'une population animale dans la nature. Imaginons un biologiste qui souhaite estimer le nombre de poissons dans un étang. Pour ce faire, il va capturer un échantillon de poissons, les marquer (par exemple, en leur attachant une petite étiquette) puis les relâcher dans l'étang. Après un certain laps de temps, il va revenir et capturer à nouveau un échantillon de poissons. En examinant combien de poissons marqués sont présents dans le deuxième échantillon, le biologiste peut estimer la taille totale de la population de poissons dans l'étang.

Soit N la taille inconnue de la population totale que nous cherchons à estimer.

1. Lors de la première capture, un échantillon de n_1 poissons est prélevé dans la population, qui sont tous marqués par le biologiste.
2. Lors de la deuxième capture, un nouvel échantillon de n_2 poissons est prélevé dans la population, le biologiste à présent ne fait que compter le nombre de poissons n_{12} qui ont été déjà marqués à la première capture.

On observe donc n_1 , n_2 et $n_{12} \leq n_2$. On cherche à estimer N . Les valeurs numériques sont: $n_1 = 125$, $n_2 = 110$ et $n_{12} = 15$.

1 Modèle à un paramètre

On ignore la probabilité de capture (le biais de sélection) et on s'intéresse uniquement à la modélisation de la recapture des poissons marqués à l'étape 2. On suppose donc que n_1 et n_2 sont des constantes données.

1. La loi hypergéométrique de paramètres (M, K, n) modélise le nombre de succès parmi n tirages sans remplacement dans une population de taille M où le nombre total de succès possible est K . Sa densité est donnée par:

$$\mathbb{P}(X = k) = \frac{\binom{K}{k} \binom{N-K}{n-k}}{\binom{N}{n}}$$

Comment peut-on l'utiliser pour modéliser $n_{12}|N$?

2. Déterminer l'estimateur de maximum de vraisemblance.
3. On se place désormais dans le cadre Bayésien. Quel problème risque de se poser en prenant une loi a priori uniforme sur \mathbb{N}^* ?
4. Prenons à présent une loi a priori uniforme donnée par $\mathbb{1}_{[n_0, N_0]}$ où N_0 est un nombre supposé assez grand. Précisez la borne inférieure n_0 .
5. Comment peut-on obtenir des statistiques a posteriori sur N numériquement ?
6. Implémentez le modèle avec `pymc`, faites le diagnostic MCMC et interprétez les résultats.

2 Modèle à deux paramètres

On modélise à présent la probabilité de capture d'un poisson par un paramètre p . n_1 et n_2 et n_{12} sont désormais des variables aléatoires. Les paramètres d'intérêt sont donc (p, N) . On pose $n'_2 = n_2 - n_{12}$.

1. Quel sont les modèles adéquats pour les variables $n_1, n_{12}|n_1$ et $n'_2|n_1$?
2. Donner la vraisemblance du modèle.
3. On suppose une loi a priori uniforme sur p, N supposés indépendants. Déterminer la distribution a posteriori de p
4. De même pour N . Pouvez-vous la reconnaître ?
5. La loi binomiale négative est une distribution de probabilité discrète définie par deux paramètres : le nombre de succès r et la probabilité de succès p . Elle modélise le nombre d'échecs indépendants nécessaires pour obtenir exactement r succès, sachant que la probabilité de succès dans chaque essai est p . Sa densité est donnée par:

$$P(X = k) = \binom{k+r}{k} \times p^r \times (1-p)^k$$

Notons $n_+ = n_1 + n'_2$. Montrez que $N - n_+|n_+, p$ suit une Binomiale négative et précisez ses paramètres.

6. Proposez une procédure pour simuler la loi a posteriori.
7. Implémentez le modèle avec `pymc`, faites le diagnostic MCMC et interprétez les résultats.
8. Comment se comparent les deux modèles ?