

Bayesian Inference with ecological applications Chapter 10

潜在的な事象を扱うための多項分布モデル *Latent Multinomial Models*

本章では、記録した頻度データが多項分布に従う潜在的な確率変数を集約したものと考えられるときの、データ解析の手法を紹介する。この手法は標識再捕獲の解析に幅広く応用でき、観察された捕獲履歴が複数の事象からなる場合には常に適切な手法である。個体の誤識別に関するモデルについて考えよう。

10.1 MODEL M_t

ある動物の閉鎖個体群における個体数 N の推定を行うことを考える。Model M_t は、その動物の捕獲確率 p_t はどの調査回 t でも個体や時間によって変化しないことを想定したモデルである。

個々の動物の捕獲履歴 capture history $\omega = (\omega_1, \omega_2, \dots, \omega_T)$ は、総調査回数 T の長さをもつ二値のベクトルにより表わされる。例えば4回の調査で1回目と3回目のみに発見できた場合 $\{1,0,1,0\}$ もしくは 1010 と記述される。捕獲履歴は

$$j \equiv j(\omega) = 1 + \sum_{t=1}^T \omega_t 2^{t-1}$$

を用いて指標化すると扱いやすい。上記の $\omega = 1010$ の例では $j = 6$ となる。

捕獲履歴 j をもつ個体の数を f_j とするとき、 $\mathbf{f} = (f_1, f_2, \dots, f_{2^T})$ は個体数 N とセル確率

$$\pi_j = \prod_{t=1}^T p_t^{\omega_{j,t}} (1-p_t)^{1-\omega_{j,t}}$$

($\omega_{j,t}$ = 捕獲履歴 j における調査回 t での記録)

を母数とする多項分布に従う確率変数である。

Model M_t による解析は容易である。観察履歴頻度 observed frequency $\mathbf{f}^+ = (f_2, f_3, \dots, f_{2^T})$ は母数に観察個体数 n とセル確率

$$\pi_j^+ = \frac{\pi_j}{1 - \pi_1}$$

をもつ多項分布に従っている。また観察個体数 n は二項分布 $n \sim B(N, 1 - \pi_1)$ に従う。多項分布に従う \mathbf{f} はそれらの積 $[\mathbf{f}] \propto [\mathbf{f}^+ | n][n]$ として表わされるだろう。PANEL 10.1 に BUGS code を示す。

10.2 MODEL $M_{t,\alpha}$

個体が正しく識別できない場合、データに存在しないのに観察される個体 ghost record が含まれるため、Model M_t は個体数を過大評価してしまう。そこで Model $M_{t,\alpha}$ では、誤識別された個体は1個体しか存在しないと想定*して、誤識別の可能性をモデルに考慮する。正識別率 correct identification probability を α とすると、Model $M_{t,\alpha}$ では以下のように3つの事象とその確率が考えられる。

$\tilde{\omega}_t = 0$	捕獲されない	確率 $1 - p_t$
$\tilde{\omega}_t = 1$	捕獲され、識別も正しい	確率 αp_t
$\tilde{\omega}_t = 2$	捕獲されたが、識別が誤っている	確率 $(1 - \alpha) p_t$

* 遺伝物質を用いた識別では妥当なようです。

したがって、このモデルでは 3^T の長さをもつ潜在履歴 latent history $\{\tilde{\omega}_1, \tilde{\omega}_2, \dots, \tilde{\omega}_T\}$ が存在し、同様に

$$i = 1 + \sum_{t=1}^T \tilde{\omega}_t 3^{t-1}$$

と指標化して扱うことができる。以下では、潜在的に起こりうる履歴を潜在履歴 latent history i 、実際にデータとして記録された履歴を記録履歴 recorded history j ということにする。

Model $M_{t,\alpha}$ の解析では、始めに潜在履歴頻度 latent history frequency $\mathbf{X} = (x_1, x_2, \dots, x_{3^T})'$ を記録履歴頻度 $\mathbf{f}^+ = (f_2, f_3, \dots, f_{2^T})$ に変換する必要がある。

$$\mathbf{f}^+ = \mathbf{A}' \mathbf{X}$$

行列 \mathbf{A} は、潜在履歴 i から記録履歴 j が生じるとき $A_{ij} = 1$ 、それ以外で $A_{ij} = 0$ となる $3^T \times (2^T - 1)$ の行列である (Table 10.1)。ここで、潜在履歴頻度の妥当な集合 $\{\mathbf{X}: \mathbf{A}'\mathbf{X} = \mathbf{f}^+\}$ を知ることができるのであれば、尤度は次式の条件付き確率の和として表現できる。

$$L(p_1, p_2, \dots, p_T, \alpha | \mathbf{f}^+) = \sum_{\mathbf{X}: \mathbf{A}'\mathbf{X} = \mathbf{f}^+} \Pr(\mathbf{X} | p_1, p_2, \dots, p_T, \alpha)$$

集合 $\{\mathbf{X}: \mathbf{A}'\mathbf{X} = \mathbf{f}^+\}$ の推定には線形代数の知識を必要とする。

Digression: Some Linear Algebra

$\mathbf{0}$ でない行列 \mathbf{B} を $\mathbf{0}$ に写像するベクトルの集合、零空間 null space という概念を導入しよう。ベクトル \mathbf{v} が行列 \mathbf{B} の零空間に属する $\mathbf{B}\mathbf{v} = \mathbf{0}, \mathbf{v} \in \text{null}(\mathbf{B})$ とき、 $\mathbf{B}\mathbf{x} = \mathbf{B}\mathbf{y}$ から次のように記述できる。

$$\mathbf{B}\mathbf{x} = \mathbf{B}\mathbf{y} \Leftrightarrow \mathbf{B}(\mathbf{x} - \mathbf{y}) = \mathbf{0} \Leftrightarrow \mathbf{x} = \mathbf{y} + \mathbf{v}$$

同様に、もし $\mathbf{A}'\mathbf{X}_0 = \mathbf{f}^+$ を満たす潜在履歴頻度のひとつ \mathbf{X}_0 を見つけることができるなら

$$\{\mathbf{X}: \mathbf{A}'\mathbf{X} = \mathbf{f}^+\} \Leftrightarrow \{\mathbf{X}: \mathbf{A}'\mathbf{X} = \mathbf{A}'\mathbf{X}_0\} \Leftrightarrow \{\mathbf{X}: \mathbf{X} = \mathbf{X}_0 + \mathbf{v}, \mathbf{v} \in \text{null}(\mathbf{A}')\}$$

と記述できる。 \mathbf{X}_0 となる潜在履歴頻度は、個体の識別が完全に正しい場合 ($\alpha = 1$) に容易に得られる。

行列 \mathbf{A}' の零空間は、線形に独立な $r = 3^T - 2^T + 1$ 個の基底ベクトル basis vector $\mathbf{b}_1, \mathbf{b}_2, \dots, \mathbf{b}_r$ からなり、 a_1, a_2, \dots, a_r を定数とすると

$$\mathbf{v} = \sum_{k=1}^r a_k \mathbf{b}_k, \forall \mathbf{v} \in \text{null}(\mathbf{B})$$

が成り立つ (Table 10.2)。

Sample Calculation of Null Space Basis Vector

$T = 2$ の場合について、行列 \mathbf{A}' の零空間における基底ベクトルの計算を例示する。行列 \mathbf{A}' は

$$\mathbf{A}' = \begin{bmatrix} 0 & 1 & 1 & 0 & 0 & 1 & 0 & 1 & 1 \\ 0 & 0 & 0 & 1 & 0 & 1 & 1 & 1 & 1 \\ 0 & 0 & 0 & 0 & 1 & 0 & 0 & 0 & 0 \end{bmatrix}$$

と表わされる。行列 \mathbf{A}' の零空間を $\mathbf{x} = (x_1, x_2, \dots, x_9)'$ とすると、 $\mathbf{A}'\mathbf{x} = \mathbf{0}$ は

$$x_2 + x_3 + x_6 + x_8 + x_9 = 0$$

$$x_4 + x_6 + x_7 + x_8 + x_9 = 0$$

$$x_5 = 0$$

が成り立つ。したがって、

$$\mathbf{x} = \begin{bmatrix} x_1 \\ x_2 \\ x_3 \\ x_4 \\ x_5 \\ x_6 \\ x_7 \\ x_8 \\ x_9 \end{bmatrix} = \begin{bmatrix} x_1 & & & & & & & & & \\ -x_3 - x_6 - x_8 - x_9 & x_1 & & & & & & & & \\ & x_3 & & & & & & & & \\ -x_6 - x_7 - x_8 - x_9 & & x_3 & & & & & & & \\ & & & 0 & & & & & & \\ & & & x_6 & & & & & & \\ & & & x_7 & & & & & & \\ & & & x_8 & & & & & & \\ & & & x_9 & & & & & & \end{bmatrix} = x_1 \begin{bmatrix} 1 \\ 0 \\ 0 \\ 0 \\ 0 \\ 0 \\ 0 \\ 0 \\ 0 \end{bmatrix} + x_3 \begin{bmatrix} 0 \\ -1 \\ 1 \\ 0 \\ 0 \\ 0 \\ 0 \\ 0 \\ 0 \end{bmatrix} + x_6 \begin{bmatrix} 0 \\ -1 \\ 0 \\ 0 \\ 1 \\ 0 \\ 0 \\ 0 \\ 0 \end{bmatrix} + x_7 \begin{bmatrix} 0 \\ 0 \\ 0 \\ -1 \\ 0 \\ 1 \\ 0 \\ 0 \\ 0 \end{bmatrix} + x_8 \begin{bmatrix} 0 \\ -1 \\ 0 \\ -1 \\ 0 \\ 0 \\ 0 \\ 1 \\ 0 \end{bmatrix} + x_9 \begin{bmatrix} 0 \\ -1 \\ 0 \\ -1 \\ 0 \\ 0 \\ 0 \\ 0 \\ 1 \end{bmatrix}$$

となり、行列 \mathbf{A}' の零空間における $6 (= 3^2 - 2^2 + 1)$ 個の基底ベクトルが求められた。

10.3 GIBBS SAMPLING FOR MODEL $M_{t,\alpha}$

個体が潜在履歴 ω_i をもつ確率：

$$\pi_i = \prod_{t=1}^T p_t^{\mathbf{I}(\omega_{i,t} > 0)} (1 - p_t)^{\mathbf{I}(\omega_{i,t} = 0)} \alpha^{\mathbf{I}(\omega_{i,t} = 1)} (1 - \alpha)^{\mathbf{I}(\omega_{i,t} = 2)}$$

個体数 N において、潜在履歴頻度 $\mathbf{X} = (x_1, x_2, \dots, x_{3T})'$ が起こる確率：

$$[\mathbf{X} | N, p_1, p_2, \dots, p_T, \alpha] = \left\{ \frac{N!}{\prod x_i!} \prod \pi_i^{x_i} \right\} \mathbf{I}(\sum x_i = N)$$

潜在履歴頻度を考慮したときの、記録履歴頻度 \mathbf{f}^+ が起こる確率：

$$[\mathbf{f}^+ | \mathbf{X}, N, p_1, p_2, \dots, p_T, \alpha] = \mathbf{I}(\mathbf{f}^+ = \mathbf{A}' \mathbf{X})$$

全条件付き分布は、事前分布を考慮すると

$$[\mathbf{f}^+ | \mathbf{X}, N, p_1, p_2, \dots, p_T, \alpha] [\mathbf{X} | N, p_1, p_2, \dots, p_T, \alpha] [N] [p_1] [p_2] \cdots [p_T] [\alpha]$$

に比例する分布である。

Prior and Full Conditionals for p_t and α

p_t と α の事前分布がベータ分布 $p_t \sim \text{Be}(a_0^t, b_0^t)$ および $\alpha \sim \text{Be}(a_0^\alpha, b_0^\alpha)$ に従うとき、全条件付き分布もまたベータ分布 $p_t \sim \text{Be}(a^t, b^t)$ および $\alpha \sim \text{Be}(a^\alpha, b^\alpha)$ に従う。

$$a^t = a_0^t + \sum_i x_i \mathbf{I}(\omega_{i,t} > 0)$$

$$b^t = b_0^t + \sum_i x_i \mathbf{I}(\omega_{i,t} = 0)$$

$$a^\alpha = a_0^\alpha + \sum_i \sum_t x_i \mathbf{I}(\omega_{i,t} > 0)$$

$$b^\alpha = b_0^\alpha + \sum_i \sum_t x_i \mathbf{I}(\omega_{i,t} = 0)$$

Prior and Full Conditionals for \mathbf{X} and N

\mathbf{X} を固定すると新しい N の値が標本抽出できないため、 (\mathbf{X}, N) の一組として Gibbs sampling を行う。この連結全条件付分布は、 $g(N)$ を N の事前分布とすると

$$[\mathbf{X}, N | \cdot] \propto \left\{ \frac{N!}{\prod x_i!} \prod \pi_i^{x_i} \right\} \mathbf{I}(\mathbf{f}^+ = \mathbf{A}' \mathbf{X}) \mathbf{I}(\sum x_i = N) g(N)$$

と表現される。 N に関する事前知識がない場合、適当な大きな値 M を与えて、事前分布に離散的な一様分布 $\{1, 2, \dots, M\}$ を用いるのが妥当かもしれない。有限な範囲をもたない非正則な一様分布 $[N] \propto 1/N$ も事後分布は正則となるので、事前分布として用いることができる。そこで、複数の事前分布を使って解析結果の感度を評価することが勧められる。

Gibbs Sampling

Gibbs sampling は次のように行われる。

- Step 1: 適当な潜在履歴頻度の集合から初期値 \mathbf{X}_{curr} を発生させる。
- Step 2: \mathbf{X}_{curr} から a^α, b^α と a^t, b^t (for $t = 1, 2, \dots, T$) を計算する。
- Step 3: p_t と α を標本抽出する。
- Step 4: $k = 0$ とおく。
- Step 5: $k = k + 1$ とする。0を除いた離散的な一様分布 $\{-D_k, \dots, -1, 1, \dots, D_k\}$ から c_k を標本抽出し、潜在履歴頻度の候補 $\mathbf{X}_{\text{cand}} = \mathbf{X}_{\text{curr}} + c_k \mathbf{b}_k$ を発生させる。ここで \mathbf{b}_k は $\text{null}(\mathbf{A})$ の基底ベクトル、 D_k は調節パラメータである。
- Step 6: $r = \min \left\{ \frac{[\mathbf{X}_{\text{cand}} | \cdot]}{[\mathbf{X}_{\text{curr}} | \cdot]}, 1 \right\}$ を計算し、 r の確率で候補値を採択する (i.e. $\mathbf{X}_{\text{curr}} = \mathbf{X}_{\text{cand}}$)。
- Step 7: $3^T - 2^T + 1$ 個の基底ベクトルすべてで、Step 5 ~ Step 6 を繰り返す。
- Step 8: Step 2 ~ Step 7 を十分に多く繰り返す。

10.4 AN IMPLEMENTATION OF MODEL $M_{t,\alpha}$

調査回数 $T = 5$ 、個体数 $N = 400$ 、捕獲確率 $p_1 = 0.3, p_2 = 0.4, p_3 = 0.5, p_4 = 0.6, p_5 = 0.7$ 、正識別率 $\alpha = 0.9$ から記録履歴頻度を発生させ、シミュレーションによる Model $M_{t,\alpha}$ の個体数推定を行った。

まず、Model M_t で平坦な事前分布を用いて推定を行ったところ、個体数 N は 25% ほどの過大評価がされた (事後分布の中央値 = 511, 95% HPDI = [495, 525])。Model $M_{t,\alpha}$ による推定は 110000 回の標本抽出と最初の 10000 回を焼き捨てにより行われた。最初の 5000 回は D_k の調節期間としても用いた。 $\delta_k = 1$ から始め、 k^{th} の基底ベクトルの新しい候補が採択された場合には δ_k を 0.95 倍し、採択されなかった場合には $1/0.95$ 倍する[†]。Step 5 の D_k には、 δ_k より大きい整数の中で最も近いものを当てはめた。

Model $M_{t,\alpha}$ の結果は満足のいくものであった (Table 10.3)。

Markov 連鎖はかなり長い自己相関を示しており (Figure 10.1)、これは Metropolis-Hastings の標本抽出における低い移動率 (標本抽出された値の変化が小さいということ?) に原因があるようだ。調節された 212 個の D_k の多くが 3 以下の値を示し、1 個 (頻度 0 の潜在履歴を増やす基底ベクトル) の D_k だけが 12 となった。さらに、この低い移動率は潜在履歴の多くが低い頻度をもつことによるようだ。243 種類の潜在履歴のうち、145 種類が頻度 0 となり、215 種類が頻度 3 以下となった。残り 28 種類の潜在履歴に注目すると実際のものとよく一致していた (Figure 10.2)。

[†] つまり、基底ベクトルの新しい候補値が採択される時は十分大きな事前分布から候補地が選ばれているとして δ_k の値を小さく、逆に基底ベクトルの新しい候補値が採択されないときは事前分布の大きさが十分でないとして δ_k の値を大きくしている。

10.5 EXTENSIONS

標識再捕獲データの多くはベルヌーイ試行に従う事象の履歴として記述される。事象の履歴が完全に観察される場合、モデルのパラメータ推定は容易だ。しかし、実際のデータの多くは、事象の履歴というよりは、それらが複合した捕獲履歴であり、その要素となる事象の多くが直接観察されない。こうした複数の事象からなる履歴は、容易な推定を許さない尤度関数を導く。本章のモデル手法は、そのような複合的な履歴が観察される多くの場合に適用できる。多項分布に従う潜在的な確率変数 \mathbf{X} を観察できない事象の履歴を記述し、変換により頻度ベクトル \mathbf{f} に集約されたものとして考えることができる。また、本章で発達させた Gibbs sampling の方法も他の標識再捕獲モデルの多くに応用できるだろう。

Model $M_{t,\alpha}$ とその他のモデルの違いは、行列 \mathbf{A} の行に含まれる 1 の数が複数になることだ。ひとつの潜在履歴頻度が複数の記録履歴頻度を生じさせる。それゆえ Model $M_{t,\alpha}$ を解析するうえで、潜在的な事象を扱う多項分布の構造に対しての知識が必要となる。こうした理解は、必ずしも他のモデルでは必要にならないが、解析に対して明瞭な概念や多目的な枠組みを与えてくれるだろう。

PANEL 10.1

```

model {
  for ( t in 1:T ) { p[ t ] ~ dunif ( 0,1 ) }           # 捕獲確率  $P_t$  の事前分布
  for ( j in 1:Cells ) {
    for ( t in 1:T ) {                                 # セル確率の計算の一部  $p_t^{\omega_{jt}} (1-p_t)^{1-\omega_{jt}}$ 
      c[ j,t ] <- pow( p[ t ] , omegas[ j,t ] ) *
      pow( 1 - p[ t ] , 1 - omegas[ j,t ] )
    }
    pi[ j ] <- prod( c[ j,1:T ] )                       # セル確率  $\pi_j : p_t^{\omega_{jt}} (1-p_t)^{1-\omega_{jt}}$  の積の計算
  }
  for ( j in 2:Cells ) {
    pi.obs[ j ] <- pi[ j ] / ( 1 - pi[ 1 ] )           # 観察履歴のセル確率  $\pi_j^+$ 
  }
  n <- sum( f[ 2: Cells ] )                             # 観察履歴から捕獲された個体数  $n$  を計算
  f[ 2: Cells ] ~ dmulti ( pi.obs[ 2:Cells ] , n )    # 観察履歴  $\mathbf{f}^+$  : ここで尤度を計算
  pn <- 1 - pi[ 1 ]
  n ~ dbin ( pn , N )                                  # 捕獲された個体数と捕獲確率から個体数  $N$  を推定
  cN ~ dunif ( n , 10000 )                             # 個体群に含まれる個体数  $N$  の事前分布
  N <- round ( cN )                                    # 個体数  $N$  を自然数に
}

```

