

## EM アルゴリズム (The Expectation Maximization)

一般にデータは不完全であり、それに基づいて推定しなければならない場合が多い。また、複雑な尤度関数の場合、最尤推定値を求めることは困難な場合がある。EM アルゴリズムは不完全データの問題を完全データのフレームワークで逐次的にパラメーターの最尤推定量を求めてゆく方法で、計算自体より実行し易いアルゴリズムである。ただし、局所的な極大に到達してしまう可能性がある。

完全データ :  $X = (Y, Z)$ , 不完全データ (観測データ) :  $Y$ , 欠失データ :  $Z$

$Y = t(X)$  :  $X$  から  $Y$  への射影

<EM アルゴリズムの手順>

パラメーターベクトル  $\mathcal{G}$  の下での  $X$  の分布密度 :  $f(X|\mathcal{G})$

①  $Q(\mathcal{G}|\mathcal{G}_n) = E[\log f(X|\mathcal{G})|Y, \mathcal{G}_n]$  : データ  $Y$ 、パラメーター  $\mathcal{G} = \mathcal{G}_n$  の下での分布密度

$f(X|\mathcal{G})$  の条件付期待値

②  $\mathcal{G}_{n+1} = \text{Arg Max}_{\mathcal{G}} Q(\mathcal{G}|\mathcal{G}_n)$  :  $Q(\mathcal{G}|\mathcal{G}_n)$  を最大にする  $\mathcal{G}$  を  $\mathcal{G}_{n+1}$  とする。

③  $g(Y|\mathcal{G})$  : 観測データ  $Y$  の尤度関数  $\log g(Y|\mathcal{G}_{n+1}) \geq \log g(Y|\mathcal{G}_n)$  が成り立つ

(観測データ  $Y$  の尤度が増加)

### § 1. EM アルゴリズムの証明

補助定理 1. (Jensen の不等式)

$X$  : 確率変数、 $h(x)$  を凸関数 (下に凸) とする。

このとき  $E[h(X)] \geq h(E[X])$  が成り立つ。厳密に凸な関数 ( $\frac{d^2}{dx^2} h(x) > 0$ ) においては等号は  $X = E[X]$  a.e. すなわち  $X$  が定数のときのみ成り立つ。

(証明)

$E[X] = \mu$  とする。  $h(x)$  は凸関数なので  $\frac{d^2}{dx^2} h(x) \geq 0$ 、テーラーの定理より

$$h(x) = h(\mu) + (x - \mu)h'(\mu) + \frac{1}{2}(x - \mu)^2 h''(\theta), \quad (\theta \in (\mu, x))$$

$$\geq h(\mu) + (x - \mu)h'(\mu)$$

$x$  を確率変数  $X$  に置き換えると  $h(W) \geq h(\mu) + (X - \mu)h'(\mu)$ 、期待値を取ると  $E[h(X)] \geq h(\mu) + (E[X] - \mu)h'(\mu) = h(\mu)$ 、すなわち  $E[h(X)] \geq h(E[X])$  が成り立つ。

$\frac{d^2}{dx^2}h(x) > 0$  のとき恒等的に  $X \equiv \mu$  でなければ、明らかに不等号  $E[h(X)] > h(E[X])$  となる。

補助定理 2. (エントロピー不等式)

$f(x), g(x)$  を測度  $\mu(x)$  に関する確率分布密度、すなわち  $\int f(x)d\mu = \int g(x)d\mu = 1$ 、

かつ  $f(x) > 0, g(x) > 0$  とする。また  $E_f[h(x)] = \int h(x)f(x)d\mu$  と定義する。

このとき  $E_f[\log f(x)] \geq E_f[\log g(x)]$ 、等号は  $f(x) = g(x)$  (a.e.  $\mu$ ) のときに限る。

(証明)  $-\log x$  は下に凸な関数なので、Jensen の不等式より

$$\begin{aligned} E_f[\log f(x)] - E_f[\log g(x)] &= E_f\left[-\log \frac{g(x)}{f(x)}\right] \geq -\log E_f\left[\frac{g(x)}{f(x)}\right] \\ &= -\log \int \frac{g(x)}{f(x)} f(x)d\mu = -\log \int g(x)d\mu = 0 \end{aligned}$$

よって  $E_f[\log f(x)] \geq E_f[\log g(x)]$ 。等号は  $\frac{g(x)}{f(x)} = E_f\left[\frac{g(x)}{f(x)}\right] = \int g(x)d\mu = 1$  が a.e.  $\mu$  に成り立つとき、すなわち  $f(x) = g(x)$  (a.e.  $\mu$ ) のときに限る。 (終り)

$$Q(\mathcal{G}|\mathcal{G}_n) = E\left[\log f(X|\mathcal{G})|Y = y, \mathcal{G}_n\right] \quad \mathcal{G}_n \text{ は現在の } \mathcal{G} \text{ の推定値}$$

$g(Y|\mathcal{G})$  : 観測データの尤度関数(likelihood)。

補助定理 3.

$$Q(\mathcal{G}_n|\mathcal{G}_n) - \log g(y|\mathcal{G}_n) \geq Q(\mathcal{G}|\mathcal{G}_n) - \log g(y|\mathcal{G})$$

(証明)  $\frac{f(x|\mathcal{G})}{g(y|\mathcal{G})} = \frac{P(X, \mathcal{G})}{P(Y, \mathcal{G})} = P(X, \mathcal{G}|Y, \mathcal{G})$  : 不完全データ  $Y$  の下での完全データ  $X$  の分布

同様に、 $\frac{f(x|\mathcal{G}_n)}{g(y|\mathcal{G}_n)} = \frac{P(X, \mathcal{G}_n)}{P(Y, \mathcal{G}_n)} = P(X, \mathcal{G}_n|Y, \mathcal{G}_n)$ 。

$A = \{x; t(x) = y\} = \{X = (y, Z)\}$  すなわち  $A$  は空間  $X = (Y, Z)$  の  $Y = y$  断面とすると、

$$\int_A \frac{f(x|\mathcal{G})}{g(y|\mathcal{G})} d\mu = \frac{1}{g(y|\mathcal{G})} \int_A f(x|\mathcal{G}) d\mu = \frac{g(y|\mathcal{G})}{g(y|\mathcal{G})} = 1, \quad \int_A \frac{f(x|\mathcal{G}_n)}{g(y|\mathcal{G}_n)} d\mu = 1$$

A 上で  $F = \frac{f(x|\mathcal{G}_n)}{g(y|\mathcal{G}_n)}$ ,  $G = \frac{f(x|\mathcal{G})}{g(x|\mathcal{G})}$  とすると、 $Fd\mu, Gd\mu$  は A 上の確率密度である。

補助定理 2 より  $E_F[\log F] \geq E_G[\log G]$ 、すなわち

$$\int_A \left( \log \frac{f(x|\mathcal{G}_n)}{g(x|\mathcal{G}_n)} \right) Fd\mu \geq \int_A \left( \log \frac{f(x|\mathcal{G})}{g(x|\mathcal{G})} \right) Gd\mu。 \text{これより}$$

$$E[\log f(x|\mathcal{G}_n) | Y = y, \mathcal{G}_n] - \log g(y|\mathcal{G}_n) \geq E[\log f(x|\mathcal{G}) | Y = y, \mathcal{G}] - \log g(y|\mathcal{G})$$

故に  $Q(\mathcal{G}_n|\mathcal{G}_n) - \log g(y|\mathcal{G}_n) \geq Q(\mathcal{G}|\mathcal{G}_n) - \log g(y|\mathcal{G})$  が成り立つ。 (証明終り)

これより、任意の  $\mathcal{G}$  に対して  $\log g(y|\mathcal{G}) - \log g(y|\mathcal{G}_n) \geq Q(\mathcal{G}|\mathcal{G}_n) - Q(\mathcal{G}_n|\mathcal{G}_n)$ 、

$\mathcal{G}_{n+1} = \text{Arg Max}_{\mathcal{G}} Q(\mathcal{G}|\mathcal{G}_n)$  とすると、 $Q(\mathcal{G}_{n+1}|\mathcal{G}_n) \geq Q(\mathcal{G}_n|\mathcal{G}_n)$  なので、

$\log g(y|\mathcal{G}_{n+1}) - \log g(y|\mathcal{G}_n) \geq Q(\mathcal{G}_{n+1}|\mathcal{G}_n) - Q(\mathcal{G}_n|\mathcal{G}_n) \geq 0$ 。故に  $\log g(Y|\mathcal{G}_{n+1}) \geq \log g(Y|\mathcal{G}_n)$

が成り立ち、EM アルゴリズムの手順②、③のステップが証明された。

<補足説明> この証明で、不等式  $\int_A \left( \log \frac{f(x|\mathcal{G}_n)}{g(x|\mathcal{G}_n)} \right) Fd\mu \geq \int_A \left( \log \frac{f(x|\mathcal{G})}{g(x|\mathcal{G})} \right) Gd\mu$  が本質的な部

分である。両辺は領域 A 上の積分であり、積分結果は  $(\mathcal{G}_n, y)$  および  $(\mathcal{G}, y)$  の関数である。

## § 2. 例

<例 1> くじ引きの当りの確率

くじ引きで、くじを最大 2 回引くチャンスがあるとする。ただし 1 回目に当たらない場合はくじを戻してもう一度引くことにする。N 人がこの方法でくじを引いて k 人が当たったとき、1 回くじを引いたとき当る確率  $\mathcal{G}$  を最尤推定する。各人の結果を当たりは A、はずれは O で表す。

不完全データ Y:  $N(A) = k, N(O) = n - k$       2 回のチャンスで当たりの人数とはずれの人数

完全データ：  $n(A), n(OA), n(OO) = N(O)$  1回目に当たった人数、2回目に当たった人数  
2回ともはずれた人数

パラメーター： 当る確率  $\mathcal{G}$

完全データ  $X$  の尤度関数

$$f(X|\mathcal{G}) = \frac{n!}{n(A)!n(OA)!n(OO)!} \mathcal{G}^{n(A)} \{(1-\mathcal{G})\mathcal{G}\}^{n(OA)} \{(1-\mathcal{G})^2\}^{n(OO)} \quad (\text{多項分布})$$

$\mathcal{G}_m$  : 第 $m$ ステップ目の  $\mathcal{G}$  の推定値

$$\begin{aligned} Q(\mathcal{G}|\mathcal{G}_m) &= E[\log f(X|\mathcal{G})|Y, \mathcal{G}_m] \\ &= E[n(A)\log \mathcal{G} + n(OA)\log((1-\mathcal{G})\mathcal{G}) + 2n(OO)\log(1-\mathcal{G}) \\ &\quad + \log\left(\frac{n!}{n(A)!n(OA)!n(OO)!}\right)|Y, \mathcal{G}_m] \end{aligned}$$

$(Y, \mathcal{G}_m)$  が与えられた条件下で  $n(OO) = N(O)$  より  $E[n(OO)|Y, \mathcal{G}_m] = N(O)$ 。

$f(X|\mathcal{G})$  は多項分布で、 $n(A) + n(OA) = N(A)$  より、 $(Y, \mathcal{G}_m)$  の条件下で  $n(A)$  は二項分布

$B(N(A), \frac{\mathcal{G}_m}{\mathcal{G}_m + (1-\mathcal{G}_m)\mathcal{G}_m})$  に従う。よって

$$\begin{aligned} n_m(A) &= E[n(A)|Y, \mathcal{G}_m] = N(A) \times \frac{\mathcal{G}_m}{\mathcal{G}_m + (1-\mathcal{G}_m)\mathcal{G}_m} = \frac{k\mathcal{G}_m}{\mathcal{G}_m + (1-\mathcal{G}_m)\mathcal{G}_m} \\ n_m(OA) &= E[n(OA)|Y, \mathcal{G}_m] = N(A) \times \frac{(1-\mathcal{G}_m)\mathcal{G}_m}{\mathcal{G}_m + (1-\mathcal{G}_m)\mathcal{G}_m} = \frac{k(1-\mathcal{G}_m)\mathcal{G}_m}{\mathcal{G}_m + (1-\mathcal{G}_m)\mathcal{G}_m} \end{aligned}$$

以上より

$$\begin{aligned} Q(\mathcal{G}|\mathcal{G}_m) &= n_m(A)\log \mathcal{G} + n_m(OA)\log((1-\mathcal{G})\mathcal{G}) + 2N(O)\log(1-\mathcal{G}) \\ &\quad + E\left[\log\left(\frac{n!}{n(A)!n(OA)!n(OO)!}\right)|Y, \mathcal{G}_m\right] \end{aligned}$$

最後の項  $E[\log(\dots)|Y, \mathcal{G}_m]$  にはパラメーター  $\mathcal{G}$  は含まれないことに注意しよう。

$Q(\mathcal{G}|\mathcal{G}_m)$  が最大になる  $\mathcal{G}$  を求める。

$$\begin{aligned} \frac{\partial Q}{\partial \mathcal{G}} &= \frac{n_m(A)}{\mathcal{G}} + n_m(OA)\left(-\frac{1}{1-\mathcal{G}} + \frac{1}{\mathcal{G}}\right) - 2N(O)\frac{1}{1-\mathcal{G}} = 0 \text{ より} \\ n_m(A) + n_m(OA) - \{n_m(A) + 2n_m(OA) + 2N(O)\}\mathcal{G} &= 0 \end{aligned}$$

$$\text{これから } \mathcal{G}_{m+1} = \frac{n_m(A) + n_m(OA)}{n_m(A) + 2n_m(OA) + 2N(O)} = \frac{N(A)}{2n - n_m(A)}$$

分母の  $n_m(A) + 2n_m(OA) + 2N(O)$  は  $n$  人が引いたくじの総数の  $m$  ステップ目の推定値であり、分子は引いた当たりくじの数である。  $n_m(A)$  を代入すると次の漸化式を得る

$$\mathcal{G}_{m+1} = \frac{k(2 - \mathcal{G}_m)}{(4n - k) - 2n\mathcal{G}_m}。 \text{極限值 } \lim_{m \rightarrow \infty} \mathcal{G}_m = \mathcal{G} \text{ とすると、 } \mathcal{G} = \frac{k(2 - \mathcal{G})}{(4n - k) - 2n\mathcal{G}}。 0 \leq \mathcal{G} \leq 1 \text{ よ}$$

$$\text{り、 } \mathcal{G} = 1 - \sqrt{1 - \frac{k}{n}} \text{ となる。}$$

(注意) この問題では2回のチャンスで当たる確率は  $\mathcal{G} + (1 - \mathcal{G})\mathcal{G}$ 、2回ともはずれの確率は  $(1 - \mathcal{G})^2$  なので、データの尤度関数は  $g(Y|\mathcal{G}) = {}_n C_k (\mathcal{G} + (1 - \mathcal{G})\mathcal{G})^k ((1 - \mathcal{G})^2)^{n-k}$  となる。

$\log g(Y|\mathcal{G}) = \log {}_n C_k + k \log \mathcal{G} + k \log(2 - \mathcal{G}) + 2(n - k) \log(1 - \mathcal{G})$ 、極値は

$$\frac{\partial}{\partial \mathcal{G}} \log g(Y|\mathcal{G}) = \frac{k}{\mathcal{G}} + \frac{k}{\mathcal{G} - 2} + \frac{2(n - k)}{\mathcal{G} - 1} = 0 \text{ より } n\mathcal{G}^2 - 2n\mathcal{G} + k = 0, (0 \leq \mathcal{G} \leq 1)。 \text{故に}$$

$\mathcal{G} = 1 - \sqrt{1 - \frac{k}{n}}$  と最尤推定値が直接もとまる。このような簡単な尤度関数の場合は直接求めることができるが、複雑な尤度関数の場合にも、逐次的に最尤推定値に近づけて行くのが EM アルゴリズムである。

## <例2> 遺伝子頻度の推定

仮定：Hardy-Weinberg の法則は成立していると仮定する。

ABO 血液型

不完全（観測）データ  $Y$ ：各血液型のサンプル数  $n(A), n(B), n(AB), n(O)$

$$n(A) + n(B) + n(AB) + n(O) = n$$

完全データ  $X$ ：遺伝子型のサンプル数  $n(AA), n(AO), n(BB), n(BO), n(AB), n(OO)$

$$n(AA) + n(AO) + n(BB) + n(BO) + n(AB) + n(OO) = n$$

パラメーター：遺伝子頻度  $\mathcal{G} = (p(A), p(B), p(O))$ ,  $p(A) + p(B) + p(O) = 1$

完全データ  $X$  の尤度関数(likelihood)

$$f(X|\mathcal{G}) = \frac{n!}{n(AA)!n(AO)!n(BB)!n(BO)!n(AB)!n(OO)!} (p(A)^2)^{n(AA)} (p(B)^2)^{n(BB)} (p(O)^2)^{n(OO)} \\ \times (2p(A)p(O))^{n(AO)} (2p(B)p(O))^{n(BO)} (2p(A)p(B))^{n(AB)}$$

$\mathcal{G}_m = (p_m(A), p_m(B), p_m(O))$ ：第  $m$  ステップでのパラメーターの値

$$\begin{aligned}
Q(\mathcal{G}|\mathcal{G}_m) &= E\left[\log f(X|\mathcal{G})|Y = (n(A), n(B), n(AB), n(O)), \mathcal{G}_m\right] \\
&= E\left[2n(AA)\log p(A) + 2n(BB)\log p(B) + 2n(OO)\log p(O) + n(AO)\log(2p(A)p(O))\right. \\
&\quad \left.+ n(BO)\log(2p(B)p(O)) + n(AB)\log(2p(A)p(B)) + \log\left(\frac{n!}{n(AA)!\cdots n(OO)!}\right) \mid Y, \mathcal{G}_m\right]
\end{aligned}$$

AB型とO型の血液型の数は観測データ Y から明らかに

$$n_m(AB) = E[n(AB)|Y, \mathcal{G}_m] = n(AB), \quad n_m(OO) = E[n(OO)|Y, \mathcal{G}_m] = n(O)$$

完全データ X の尤度は多項分布なので、Hardy-Weinberg の法則は成り立つという仮定の下で、条件 Y,  $\mathcal{G}_m$  の下で、遺伝子型 AA および AO の数の条件付期待値は

$$\begin{aligned}
n_m(AA) &= E[n(AA)|Y, \mathcal{G}_m] = n(A) \times \frac{\text{Freq of AA}}{\text{Freq of A type}} = n(A) \times \frac{(p_m(A))^2}{(p_m(A))^2 + 2p_m(A)p_m(O)} \\
n_m(AO) &= E[n(AO)|Y, \mathcal{G}_m] = n(A) \times \frac{2p_m(A)p_m(O)}{(p_m(A))^2 + 2p_m(A)p_m(O)}
\end{aligned}$$

A を B に置き換えれば同様に  $n_m(BB), n_m(BO)$  も求められる。

$$\begin{aligned}
Q(\mathcal{G}|\mathcal{G}_m) &= 2n_m(AA)\log p(A) + 2n_m(BB)\log p(B) + 2n_m(OO)\log p(O) \\
&\quad + n_m(AO)(\log 2 + \log p(A) + \log p(O)) + n_m(BO)(\log 2 + \log p(B) + \log p(O)) \\
&\quad + n_m(AB)(\log 2 + \log p(A) + \log p(B)) + E\left[\log\left(\frac{n!}{n(AA)!\cdots n(OO)!}\right) \mid Y, \mathcal{G}_m\right]
\end{aligned}$$

最後の項の期待値はパラメーター  $\mathcal{G}$  を含まないので次のステップの  $\mathcal{G}$  の推定に関係しない。

$Q(\mathcal{G}|\mathcal{G}_m)$  を最大にする  $\mathcal{G}$  の値が  $\mathcal{G}_{m+1}$  となる。ラグランジュ未定係数法によって求める。

$$H(\mathcal{G}, \lambda) = Q(\mathcal{G}|\mathcal{G}_m) + \lambda(p(A) + p(B) + p(O) - 1) \quad \text{とする。}$$

$$\frac{\partial}{\partial p(A)} H = \frac{2n_m(AA)}{p(A)} + \frac{n_m(AO)}{p(A)} + \frac{n_m(AB)}{p(A)} + \lambda = 0 \quad \dots(1)$$

$$\frac{\partial}{\partial p(B)} H = \frac{2n_m(BB)}{p(B)} + \frac{n_m(BO)}{p(B)} + \frac{n_m(AB)}{p(B)} + \lambda = 0 \quad \dots(2)$$

$$\frac{\partial}{\partial p(O)} H = \frac{2n_m(OO)}{p(O)} + \frac{n_m(AO)}{p(O)} + \frac{n_m(BO)}{p(O)} + \lambda = 0 \quad \dots(3)$$

$$\frac{\partial}{\partial \lambda} H = p(A) + p(B) + p(O) - 1 = 0 \quad \dots(4)$$

(1),(2),(3)より

$$p(A) = (2n_m(AA) + n_m(AO) + n_m(AB)) / (-\lambda), \quad p(B) = (2n_m(BB) + n_m(BO) + n_m(AB)) / (-\lambda)$$

$$p(O) = (2n_m(OO) + n_m(AO) + n_m(BO)) / (-\lambda)$$

(4)に代入して  $\lambda = -2(n_m(OO) + n_m(AA) + \dots) = -2n$ 。

これより

$$p_{m+1}(A) = \frac{2n_m(AA) + n_m(AO) + n_m(AB)}{2n}, \quad p_{m+1}(B) = \frac{2n_m(BB) + n_m(BO) + n_m(AB)}{2n}$$

$$p_{m+1}(O) = \frac{2n_m(OO) + n_m(AO) + n_m(BO)}{2n} \quad \text{Gene counting}$$

<まとめ>

血液型 観測データ :  $Y = \{n(A), n(B), n(AB), n(O)\}$

パラメータ :  $\mathcal{G} = \{p(A), p(B), p(O)\}$  遺伝子頻度

初期推定値 :  $\mathcal{G}(0) = \{p_0(A), p_0(B), p_0(O)\}$  任意の推定値

第mステップ :  $\mathcal{G}_m = \{p_m(A), p_m(B), p_m(O)\}$

$$n_m(AA) = n(A) \times \frac{(p_m(A))^2}{(p_m(A))^2 + 2p_m(A)p_m(O)},$$

$$n_m(AO) = n(A) \times \frac{2p_m(A)p_m(O)}{(p_m(A))^2 + 2p_m(A)p_m(O)}$$

などにより各遺伝子型の数を推定

第m+1ステップ :  $\mathcal{G}_{m+1} = \{p_{m+1}(A), p_{m+1}(B), p_{m+1}(O)\}$

$$p_{m+1}(A) = \frac{2n_m(AA) + n_m(AO) + n_m(AB)}{2n} \quad \text{など Gene counting により}$$

$p_{m+1}(B), p_{m+1}(O)$  も求める。

<例3> Haplotype 頻度の推定

二つの遺伝子座 : 対立遺伝子 遺伝子座1  $\{A, a\}$ 、 遺伝子座2  $\{B, b\}$

Hardy-Weinberg 平衡を仮定、連鎖平衡は仮定しない。

Haplotype 頻度 :  $P(AB), P(Ab), P(aB), P(ab)$  を推定する。

サンプル個体数 =  $n$

不完全 (観測) データ  $Y$  :  $n(AABB), n(AABb), n(AAbb), n(AaBB), n(AaBb),$   
 $n(Aabb), n(aaBB), n(aaBb), n(aabb)$

完全データ  $X$  :  $n\left(\frac{AB}{AB}\right), n\left(\frac{AB}{Ab}\right), n\left(\frac{AB}{aB}\right), n\left(\frac{AB}{ab}\right), n\left(\frac{Ab}{aB}\right), n\left(\frac{Ab}{Ab}\right), n\left(\frac{Ab}{ab}\right), n\left(\frac{aB}{aB}\right),$   
 $n\left(\frac{aB}{Ab}\right), n\left(\frac{aB}{ab}\right), n\left(\frac{ab}{ab}\right)$

不完全データでは  $n(AaBb)$  と表示される両遺伝子座ヘテロ接合体の個体が完全データでは  $n\left(\frac{Ab}{aB}\right), n\left(\frac{AB}{ab}\right)$  と相が区別されている。

パラメーター  $\mathcal{G} = (P(AB), P(Ab), P(aB), P(ab))$

完全データの尤度関数(likelihood) 多項分布により

$$f(X|\mathcal{G}) = \frac{n!}{n\left(\frac{AB}{AB}\right)! \cdots n\left(\frac{ab}{ab}\right)!} (P^2(AB))^{n\left(\frac{AB}{AB}\right)} (2P(AB)P(Ab))^{n\left(\frac{AB}{Ab}\right)} \cdots (P^2(ab))^{n\left(\frac{ab}{ab}\right)}$$

$$Q(\mathcal{G}|\mathcal{G}_m) = E\left[\log f(X|\mathcal{G}) | Y, \mathcal{G}_m\right] = E\left[ \begin{array}{l} 2n\left(\frac{AB}{AB}\right) \log P(AB) + n\left(\frac{AB}{Ab}\right) \log(2P(AB)P(Ab)) + \\ \cdots + 2n\left(\frac{ab}{ab}\right) \log P(ab) + \log \frac{n!}{n\left(\frac{AB}{AB}\right)! \cdots n\left(\frac{ab}{ab}\right)!} \end{array} \middle| Y, \mathcal{G}_m \right]$$

Double hetero  $n\left(\frac{AB}{ab}\right), n\left(\frac{Ab}{aB}\right)$  以外の遺伝子頻度は不完全データと完全データの違いは

ないので、例えば  $n\left(\frac{AB}{AB}\right) = n(AABB)$ ,  $n\left(\frac{AB}{Ab}\right) = n(AABb)$  など全て観測データから求まる。

故に第  $m$  ステップでの Double hetero  $n\left(\frac{AB}{ab}\right), n\left(\frac{Ab}{aB}\right)$  以外の遺伝子型の数については

$$n_m\left(\frac{AB}{AB}\right) = E\left[n\left(\frac{AB}{AB}\right) | Y, \mathcal{G}_m\right] = n(AABB) \text{ 等となる。}$$

Double hetero  $n\left(\frac{AB}{ab}\right), n\left(\frac{Ab}{aB}\right)$  については多項分布の性質より

$$n_m\left(\frac{AB}{ab}\right) = E\left[n\left(\frac{AB}{ab}\right) | Y, \mathcal{G}_m\right] = n(AaBb) \times \frac{2P_m(AB)P_m(ab)}{2P_m(AB)P_m(ab) + 2P_m(Ab)P_m(aB)}$$

$$n_m\left(\frac{Ab}{aB}\right) = E\left[n\left(\frac{Ab}{aB}\right) | Y, \mathcal{G}_m\right] = n(AaBb) \times \frac{2P_m(Ab)P_m(aB)}{2P_m(AB)P_m(ab) + 2P_m(Ab)P_m(aB)}$$

例 2 と同様にしてラグランジュの未定乗数法により、各ハプロイド頻度の推定値は

$$P_{m+1}(AB) = \frac{1}{2n} \left\{ 2n(AABB) + n(AABb) + n(AaBB) + n_m\left(\frac{AB}{ab}\right) \right\}$$

$$P_{m+1}(Ab) = \frac{1}{2n} \left\{ n(AABb) + 2n(AAbb) + n_m\left(\frac{Ab}{aB}\right) + n(Aabb) \right\}$$



$$P_{m+1}(aB) = \frac{1}{2n} \left\{ n(AaBB) + n_m \left( \frac{aB}{Ab} \right) + 2n(aaBB) + n(aaBb) \right\}$$

$$P_{m+1}(ab) = \frac{1}{2n} \left\{ n_m \left( \frac{ab}{AB} \right) + n(Aabb) + n(aaBb) + 2n(aabb) \right\}$$

また A,B の遺伝子頻度は

$$P_m(A) = P_m(AB) + P_m(Ab)$$

$$= \frac{1}{2n} \{ 2n(AABB) + 2n(AABb) + 2n(AAbb) + n(AaBB) + n(AaBb) + n(Aabb) \}$$

$$P_m(B) = P_m(AB) + P_m(aB)$$

$$= \frac{1}{2n} \{ 2n(AABB) + 2n(AaBB) + 2n(aaBB) + n(AABb) + n(AaBb) + n(aaBb) \}$$

これらは観測データから決まるので m に依存しない定数である。

$P_m(AA) + P_m(Ab) + P_m(aB) + P_m(ab) = 1$  なので、3 変数を推定すれば十分である。

<参考> ラグランジュの未定乗数法

二つの  $C^1$  級関数  $F(x, y), G(x, y)$  に対して、 $H(x, y, \lambda) = F(x, y) - \lambda G(x, y)$  と置く。  
条件  $G(x, y) = 0$  の下で、関数  $F(x, y)$  は点  $(x_0, y_0)$  で極値を取るとする。さらに、

$G_x(x_0, y_0), G_y(x_0, y_0)$  の少なくとも一つが 0 でないならば、

$$H_x(x_0, y_0, \lambda) = F_x(x_0, y_0) - \lambda G_x(x_0, y_0) = 0,$$

$$H_y(x_0, y_0, \lambda) = F_y(x_0, y_0) - \lambda G_y(x_0, y_0) = 0$$

を満たす  $\lambda$  が存在する。

(証明) 条件より  $G_x(x_0, y_0), G_y(x_0, y_0)$  の少なくとも一つが 0 でないので  $G_y(x_0, y_0) \neq 0$

の場合を示す。このとき陰関数の定理より点  $x_0$  の近傍で  $C^1$  級関数  $y = f(x)$  で

$G(x, f(x)) = 0$  および  $f(x_0) = y_0$  を満たすものが存在する。  $\varphi(x) = F(x, f(x))$  とすると、

$\varphi(x)$  は点  $x_0$  で極値を取るので  $\varphi'(x_0) = 0$ 、すなわち

$$F_x(x_0, f(x_0)) + F_y(x_0, f(x_0))f'(x_0) = 0, \quad G(x, f(x)) = 0 \text{ の両辺を } x \text{ で微分すると}$$

$$G_x(x_0, f(x_0)) + G_y(x_0, f(x_0))f'(x_0) = 0。二つの式から  $f'(x_0)$  を消去すると$$

$$F_x(x_0, f(x_0)) - F_y(x_0, f(x_0)) \frac{G_x(x_0, y_0)}{G_y(x_0, y_0)} = 0. \quad \frac{F_y(x_0, y_0)}{G_y(x_0, y_0)} = \frac{F_x(x_0, y_0)}{G_x(x_0, y_0)} = \lambda \text{ と置くと、}$$

$$F_x(x_0, f(x_0)) - \lambda G_x(x_0, f(x_0)) = 0, \quad F_y(x_0, f(x_0)) - \lambda G_y(x_0, f(x_0)) = 0 \text{ が成り立つ。}$$

$\lambda$  を未定乗数と言ひ、この方法をラグランジュの未定乗数法と言ふ。(証明終わり)

例題： 条件  $x^2 + y^2 = 1$  の下で、関数  $x^2 y^3$  の最大値、最小値を求めよ。

(解)  $H(x, y, \lambda) = F(x, y) - \lambda G(x, y) = x^2 y^3 - \lambda(x^2 + y^2 - 1)$  とする。

条件  $G(x, y) = x^2 + y^2 - 1 = 0$  と  $G_x(x, y) = 2x = 0$ ,  $G_y(x, y) = 2y = 0$  を同時に満たす

$(x, y)$  はないので、

$$H_x(x_0, y_0, \lambda) = 2x_0 y_0^3 - 2\lambda x_0 = 0, \quad \text{を満たす } \lambda \text{ が存在する。}$$

$$H_y(x_0, y_0, \lambda) = 3x_0^2 y_0^2 - 2\lambda y_0 = 0$$

$x_0 \neq 0, y_0 \neq 0$  のとき、 $2\lambda = 2y_0^3 = 3x_0^2 y_0$ , 故に  $2y_0^2 = 3x_0^2$ 。  $x_0^2 + y_0^2 = 1$  なので

$$x_0^2 = \frac{2}{5}, \quad y_0^2 = \frac{3}{5} \quad \text{すなわち } (x_0, y_0) = \left( \sqrt{\frac{2}{5}}, \pm \sqrt{\frac{3}{5}} \right), \left( -\sqrt{\frac{2}{5}}, \pm \sqrt{\frac{3}{5}} \right),$$

$$\text{このとき、 } F\left(\pm \sqrt{\frac{2}{5}}, \sqrt{\frac{3}{5}}\right) = \frac{6}{25} \sqrt{\frac{3}{5}}, \quad F\left(\pm \sqrt{\frac{2}{5}}, -\sqrt{\frac{3}{5}}\right) = -\frac{6}{25} \sqrt{\frac{3}{5}} \text{。}$$

$x_0 = 0$  または  $y_0 = 0$  のとき  $F(x_0, y_0) = 0$ 、よつて

$$\text{最大値 } F\left(\pm \sqrt{\frac{2}{5}}, \sqrt{\frac{3}{5}}\right) = \frac{6}{25} \sqrt{\frac{3}{5}}, \quad \text{最小値 } F\left(\pm \sqrt{\frac{2}{5}}, -\sqrt{\frac{3}{5}}\right) = -\frac{6}{25} \sqrt{\frac{3}{5}} \text{。}$$