

## 脱落のある寿命データの数理解析

丸大 理 柳川 克

### 1. 序

1.1. 脱落のある寿命データ 生命科学における寿命データは動物実験データと臨床医学データに大別できる。両者の特長はつぎの通りである。

動物実験データ：例として発がん性が疑われている化学物質の発がん性テストを考えよう。実験動物はランダムにK群に分けられ、各群に対して相異なる量の化学物質が割りつけられ、ランダムに各個体に化学物質が投与される。厳密な管理の下で腫瘍発症までの日数が観察される。かかる状況の下では各群内の観測値は互いに独立、同一分布に従う確率変数の実現値とみなすことができ、数学モデルを導入できる。しかしながら、実験途中で実験動物が他の原因で死亡したり、ある期間を経過した後も観察が打ち切られることが多く得られるデータは不完全データである場合が多い。

臨床医学データ：ある事象（治癒、がんの転移、死亡など）が発生するまでの期間を観察し、得られるデータに基づいて治療効果が比較される。このとき、患者が突然来院しなくなる、たり、転院したりするに動物実験の場合と同様のデータの脱落、打ち切りの問題が頻繁に生じるほか、ヒトの生存には年齢、性別、生活環境などさまざまな因子が関与してくるから、これらを考慮した数学モデルを立ててデータ解析を行うことが必要がある。

これら2種類のデータは、寿命の長さを表わす変量のほか、 $p$ 個の変量  $Z = (Z_1, Z_2, \dots, Z_p)$  を導入することによって統一的に取り扱うことができる。例えば、動物実験データでは  $p=1$  とし、 $Z_1 = 1$  (群1),  $Z_1 = 2$  (群2),  $\dots$ ,  $Z_1 = K$  (群K) なるダミー変数  $Z_1$  を導入すればよい。また、臨床医学データでは、このほか、 $Z_2$  を性別、 $Z_3$  を年齢、 $\dots$  などと表わす変量に定めればよい。

2節が示すように、打ち切りは脱落の特別の場合として定式化できる。従って、本稿では脱落や打ち切りのあるデータを総称して脱落のあるデータとよぶことにする。脱落のあるデータは、少くとも自体がある期間生存したという情報を持つている。これを捨てることは、せっかく得られた情報と利用しきれずばかりか、結果にバイアスを与えることになる。

なりかわり。

1.2. 本稿の目的 Cox (1972) は脱落のある寿命データと、共変量を考慮しながら解析する強力な方法を提言した。Cox法による生命科学データの解析は英米で大流行がある。Cox法の数学的側面を明らかにする研究も盛んになっている。特に従来のノンパラメトリック法の流れの中でCox法を位置付ける試みは面白そうであるが、部分的には成功してはいない。本稿では、寿命分布の連続性がある場合に話を制限し、Cox法の数学的側面に焦点を当て、その特長を考察することを目的とする。

## 2. 定式化

脱落のある寿命データを数学的につぎのように定式化する。

•  $T_1^0, T_2^0, \dots, T_N^0$  ; 互いに独立な確率変数

$T_i^0$  は個体  $i$  の真の寿命に対応する確率変数である。脱落があるため観測可能とは限らない。

•  $F_i(t) = P\{T_i^0 \geq t\}$  ; 個体  $i$  の生存関数、連続性を仮定する。

•  $f_i(t) = -dF_i(t)/dt$

•  $\lambda_i(t) = \frac{f_i(t)}{F_i(t)} = \lim_{\Delta \rightarrow 0} \frac{P\{t \leq T_i^0 < t+\Delta | T_i^0 \geq t\}}{\Delta}$

$\lambda_i(t)$  はハザード関数とよばれる。個体  $i$  の時刻  $t$  における瞬間死亡率<sup>\*</sup>である。つぎの関係式が成立することは明らかである。

$$F_i(t) = \exp\left\{-\int_0^t \lambda_i(x) dx\right\} \quad (2.1)$$

- $U_1, U_2, \dots, U_N$  ;  $T_i^0$  に独立な確率変数,  $T_i^0$  と独立。
- $H_i(t) = P\{U_i \geq t\}$  ; 連続性を仮定可。
- $f_i(t) = -dH_i(t)/dt$

$U_i$  は個体  $i$  の脱落に用いる確率変数である。

- $T_i = \min(T_i^0, U_i)$
- $\varepsilon_i = \begin{cases} 1 & ; T_i = T_i^0 \quad (\text{死亡}) \\ 0 & ; T_i > U_i \quad (\text{脱落}) \end{cases}$

脱落のある寿命データでは  $(T_i, \varepsilon_i)$ ,  $i=1, 2, \dots, N$  が観測可能である。問題は  $(T_i, \varepsilon_i)$ ,  $i=1, 2, \dots, N$  を観測して得られたデータに基づいて  $\lambda_i(t)$  の推定,  $\lambda_i(t)$  に関する仮説の検定, およびリスク要因の分析を行なうこと、あるいは、リスク要因を分析し、それに基づいて調整し、そこから生存関数  $F_i(t)$  を推定したり、検定したりすることである。ここで、 $F_i(t)$  として、工業製品の寿命試験の類似として、未知母数を持つ指数分布、ワイブル分布、対数正規分布など、 $H_i(t)$  として指

---

\* 本稿から事象の発生を“死亡”とよぶことにする。

数分布や一様分布などが考えられる。しかし、一般に生命科学では生物の個体差、実験や観察の理想状態が行なえる<sup>429</sup>ため特定の分布型が指定できる。\$F\_i, H\_i\$ は未知とされることが多い。ここでは、\$F\_i, H\_i\$ に特定の分布型を指定し、\$H\_i\$ があればノンパラメトリックな枠組で問題を考えようというわけがある。

PROPOSITION

$$P\{T \geq t, \varepsilon_j = 1\} = \int_t^{\infty} f_j(x) H_j(x) dx$$

$$P\{T \geq t, \varepsilon_j = 0\} = F_j(t) H_j(t) - \int_t^{\infty} f_j(x) H_j(x) dx$$

$$\lim_{\Delta \rightarrow 0} \frac{P\{t \leq T < t + \Delta, \varepsilon_j = 1\}}{\Delta} = f_j(t) H_j(t)$$

$$\lim_{\Delta \rightarrow 0} \frac{P\{t \leq T < t + \Delta, \varepsilon_j = 0\}}{\Delta} = F_j(t) h_j(t)$$

(証明は簡単な省略)

Proposition 5) \$(T\_i, \varepsilon\_i), i=1, 2, \dots, N\$ の同時密度関数は次式で与えられる。

$$\begin{aligned} L_0(t_1, \dots, t_N, \varepsilon_1, \dots, \varepsilon_N) &= \prod_{j=1}^N [f_j(t_j) H_j(t_j)]^{\varepsilon_j} [F_j(t_j) h_j(t_j)]^{1-\varepsilon_j} \\ &= \left\{ \prod_{j=1}^N [\lambda_j(t_j)]^{\varepsilon_j} \exp\left[-\int_0^{t_j} \lambda_j(t) dt\right] \right\} \left\{ \prod_{j=1}^N [H_j(t_j)]^{\varepsilon_j} [h_j(t_j)]^{1-\varepsilon_j} \right\} \end{aligned}$$

よって \$(t\_i, \varepsilon\_i)\_{i=1 \sim N}\$ が与えられ \$T\_i\$ と \$\varepsilon\_i\$ の全密度関数がある。

\$H\_j\$ に依存しない推論式を得るには \$U\_i, i=1 \sim N\$ と与えられ \$T\_i\$ と \$\varepsilon\_i\$ の \$(T\_i, \varepsilon\_i), i=1 \sim N\$ の条件付き密度関数から求まる。これは簡単な計算より次式で与えられる。以後、本稿では \$T\_i\$ の \$L\_c\$ に基づいて考える。

$$L_c(t_1, \dots, t_N; \varepsilon_1, \dots, \varepsilon_N) = \prod_{j=1}^N [\lambda_j(t_j)]^{\varepsilon_j} \exp\left[-\int_0^{t_j} \lambda_j(t) dt\right] \quad (2.2)$$

### 3. Cox の方法

Cox (1972) の方法は 2, の漸近的なアイデアよりなる。  
その一つは PH モデルの導入、他は部分尤度法の導入である。

#### 3.1. PH モデル

$$\ln(\lambda_i(t)/\lambda_0(t)) = \beta_1 z_{i1} + \dots + \beta_p z_{ip} = \beta' z_i \quad (3.1)$$

と置き、相対化された瞬間死亡率の対数とその変数  $z_1, \dots, z_p$  の線形結合で説明しようというものである。Cox (1958) のロジスティックモデルの応用にはかならない。  $\lambda_0$  は  $z$  の基底に対する瞬間死亡率である。全日平均瞬間死亡率などが利用できる場合、それと  $\lambda_0$  とすれば問題は簡単であるが、ここでは  $\lambda_0$  は未知である。  $z_i$  から平均を差引いた  $z_i - \bar{z}_i$  と改めた  $z_i$  と置き、  $z_i = 0$  のときの瞬間死亡率が  $\lambda_0$  であると理解すれば便利がある。Cox は  $z$  が時間に関係する変数を含むことと許しているが、モデルの一意性をいう点で若干の問題があるようである。ここでは、  $z$  は時間に関係する項を含むものとしておく。

(3.1) を (2.2) 式に代入し対数部分の尤度関数を求めると

$$\ln L_c(t, \varepsilon) = \sum_{j=1}^N [\varepsilon_j \beta' z_j + \varepsilon_j \ln \lambda_0(t_j)] - \sum_{j=1}^N e^{\beta' z_j} \int_0^{t_j} \lambda_0(x) dx \quad (3.2)$$

$\lambda_0$  に対応する密度関数を  $f_0$  とする。  $f_0$  として、互に無関係な未知母数を含む指数分布、ワイブル分布、対数正規分布などを想定し、これらの未知母数と  $\beta$  に関して (3.2) 式を最大にするることによって  $\beta$  の推定量を求めよることが出来る。実際松原・後藤 [1979] は  $f_0$  としてベキ正規分布を仮定し、適応的に分布型を推定しながら  $\beta$  の推定を行なう意欲的な方法を提案している。しかしながら、これはつぎの2点でやっかいな問題を含んでいる。

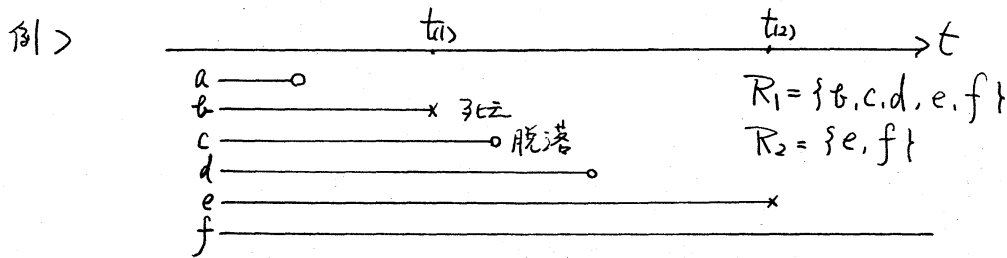
(1) 計算が (Cox法に比べて) 極めて面倒である。

(2) PHモデルでは我々の関心は、単に  $\beta$  に関する推定推定にある。このとき、 $\lambda_0$  は局外母数である。局外母数が存在するときの最尤推定量は、特にサンプル数が少ないとき、よく知られているように理論的に問題があり信頼が乏しい場合が多い。また、分布形の推定から生じる specification error の影響も無視ができない。

### 3.2 部分尤度法

まず、つぎのようにデータをいわゆる生命表形式にまとめることより始める。

(3.3)  $\left\{ \begin{array}{l} t_{(1)} < t_{(2)} < \dots < t_{(k)} ; \text{ 死亡時刻} \\ i_j ; t_{(j)} \text{ で死亡した個体名} \\ R_j = R(t_{(j)}) ; \text{ 時刻 } t_{(j)} - 0 \text{ で生存しており、かつ観察} \\ \text{下にある個体名のリスト、リスク集合} \\ \text{という。} \end{array} \right.$



注1 >  $F_N$  の連続性が仮定されるが、実際には死亡時刻は日ごとの記録しか行わないため同時刻に複数個の個体が死亡することもあり得る。このときの取り扱いは付録にゆかり、ここでは簡単のため重複がないものとしておく。

注2 > 上のようデータをとめることは、区間  $(t_{(j)}, t_{(j+1)})$  で生じたすべての脱落を時刻  $t_{(j+1)}$  で生じたと見做し、脱落に関する情報を丸めておくことに注意しよう。

(3.3) のデータ  $E$ 、 $R_j$  が所与で、集合  $R_j$  の中からランダムに1個個体  $i$  を抽出すると  $i$  が抽出された個体が  $i_j$  であると見做す。  $R_j$  の中から  $i_j$  が抽出される確率は

$$\lambda_{i_j}(t_{(j)}) / \sum_{i \in R_j} \lambda_i(t_{(j)})$$



これ  $j = 1, 2, \dots, l$  に  $\lambda_0$  の掛け合わせたのが Cox の部分尤度関数である。即ち、部分尤度関数は次式で与えられる。

$$L_P(\beta) = \prod_{j=1}^l \frac{\lambda_{ij}(t_{ij})}{\sum_{i \in R_j} \lambda_i(t_{ij})}$$

PH モデル (3.1) を代入し、対数尤度を求めると

$$\ln L_P(\beta) = \sum_{j=1}^l \beta' z_{ij} - \sum_{j=1}^l \ln \sum_{i \in R_j} e^{\beta' z_i}$$

Cox の部分尤度法は、 $\beta$  に関する推定や検定を行なう為には  $\ln L_P(\beta)$  を普通の対数尤度関数と見做し、従来の最尤法を適用せよと主張する。上式は  $\lambda_0$  を含んでいない。よって、Cox の部分尤度法では、未知のハザード関数  $\lambda_0$  に依存することなく  $\beta$  に対する統計的推測を行なうことができるわけである。

#### 4. 生存関数の推定

$\beta$  の推定値が求められたら、つぎに求められるのは共変量  $Z$  で調整したときの生存関数の推定である。それには  $\lambda_0$  をデータから推定しておかねばならない。

(2.2) の条件の下尤度  $L_C$  は、データが (3.3) の形に与えられると、注2で述べたことから注意して計算すると、つぎのように表現ができる。

$$L_c^* = \prod_{j=1}^q \lambda_{i_j}(t_{j-1}) \exp \left\{ - \sum_{j=1}^q \int_{t_{j-1}}^{t_j} \sum_{i \in R_j} \lambda_i(t) dt \right\} \quad (4.1)$$

PH モデル (3.1) の  $\lambda$  は、さきに変形すると

$$L_c^* = L_p(\beta) \cdot L_D(\beta, \lambda_0)$$

と表わされる。ここに、 $L_p$  は Cox の部分尤度関数、 $L_D$  は次式で与えられる。

$$L_D(\beta, \lambda_0) = \left[ \prod_{j=1}^q \sum_{i \in R_j} e^{\beta' z_i} \lambda_0(t_{j-1}) \right] \exp \left[ - \sum_{j=1}^q \sum_{i \in R_j} e^{\beta' z_i} \int_{t_{j-1}}^{t_j} \lambda_0(t) dt \right]$$

$L_p$  には  $\lambda_0$  は含まれていない。よって、 $L_c^*$  が  $\lambda_0$  によって最大にするには、 $L_D$  が最大にする  $\lambda_0$  を求めなければならない。本来は  $L_c^*$  が最大にする  $(\beta, \lambda_0)$  が  $(\beta, \lambda_0)$  の (条件つき) 最大推定量かあるか、ここで言うのは、 $\beta$  の推定値  $\beta^*$  が部分尤度  $L_p(\beta)$  が最大にする  $\beta$  によって求められておると考え、 $\beta = \beta^*$  が既知として  $L_D(\beta^*, \lambda_0)$  が  $\lambda_0$  によって最大にする  $\lambda_0$  の推定量を求めようと考えることである。これは Breslow (1972) のタイプ I がある。Cox (1972) も  $\lambda_0$  の推定法を提案しているが、理論的に根拠が薄弱であるように思われる。両者の比較はすでに行われている。また、 $\beta$  が主役、 $\lambda_0$  が脇役の母数があるとして Breslow 法は直観的に極めて妥当な推定法と思われ、その正当性は如何にして保証できるであろうか。  $\lambda_0$  の推定に焦点を当てて

の問題を考察することは面白いところがある。これにより、Fraser (1968) が導いた marginal likelihood の考えが役に立ちそうである。実際、Kalbfleisch and Prentice (1973) は marginal likelihood の立場から Cox の部分尤度関数と理論的に導びく試みを行っている。

とすると、

$$\lambda_0(t) = \lambda_i \quad t_{(i-1)} < t \leq t_{(i)}, \quad i=1, 2, \dots, l$$

と仮定する。このとき

$$L_D(\beta^*, \lambda_0) = \left[ \prod_{j=1}^l \sum_{i \in R_j} e^{\beta^* z_i} \lambda_j \right] \exp \left[ - \sum_{j=1}^l \sum_{i \in R_j} e^{\beta^* z_i} \lambda_j (t_j - t_{j-1}) \right]$$

$$\partial L_D(\beta^*, \lambda_0) / \partial \lambda_\alpha = 0 \quad \text{より}$$

$$\hat{\lambda}_\alpha = \left\{ [t_{(\alpha)} - t_{(\alpha-1)}] \sum_{i \in R_\alpha} e^{\beta^* z_i} \right\}^{-1}$$

を得る。これを (2.1) 式に代入すると生存関数の推定量

$$\hat{H}_0(t_{(k)}) = \exp \left\{ - \sum_{j=1}^k \frac{1}{\sum_{i \in R_j} e^{\beta^* z_i}} \right\} \doteq \prod_{j=1}^k [1 - \hat{\pi}_j]$$

$$\text{を得る。ここから、} \hat{\pi}_j = \left[ \sum_{i \in R_j} e^{\beta^* z_i} \right]^{-1}$$

特に  $\beta^* = 0$  のとき、 $\hat{H}_0$  は良く知られた Kaplan-Meier 推定量 (1958) と一致する。即ち、Breslow 法は共変量を調整した場合の Kaplan-Meier 推定量の一般化を提供する。

$F_i(t)$  の推定量は (3.1) 式より下記のようによい。

$$\hat{h}_i(t_{(k)}) = \exp\left\{-e^{\beta'Z_i} \left[ \frac{t_k}{\sum_{j=1}^k \frac{1}{e^{\beta'Z_j}}} \right]\right\}$$

### 5. 部分尤度法の正当化

実は Cox (1972) では  $L_p$  は "条件,  $\lambda$  尤度" とよばれて  
 いる。それからどのように意味で条件 $\lambda$ 尤度であるのか、 $L_p$   
 とは、この尤度関数とみれば通常の方法と適用すればいい  
 という根拠は何か、という数々の疑問が集中した。これに対  
 して Cox (1975) は "条件,  $\lambda$  尤度" という言葉を取り消し、  
 部分尤度という概念を導入して理論的正当化を行った。そ  
 の根拠はつぎの通りである。

$X_i = (R_i, t_i)$  とおく。事象は経時的に  $(X_1, i_1), (X_2, i_2),$   
 $\dots, (X_d, i_d)$  と観察されることより、(3.3) のデータを得る  
 同時確率密度関数は

$$\begin{aligned} f_{\beta, \lambda_0}(i_1, \dots, i_d; R_1, \dots, R_d; t_1, \dots, t_d) &= f_{\beta, \lambda_0}(X_1) f_{\beta, \lambda_0}(i_1 | X_1) f_{\beta, \lambda_0}(X_2 | \\ & i_1, X_1) f_{\beta, \lambda_0}(i_2 | i_1, X_1, X_2) \dots = \prod_{\alpha=1}^d f_{\beta, \lambda_0}(i_\alpha | i_1, \dots, i_{\alpha-1}; X_1, \dots, X_\alpha) \\ & \prod_{\alpha=1}^d f_{\beta, \lambda_0}(X_\alpha | i_1, \dots, i_{\alpha-1}; X_1, \dots, X_{\alpha-1}) \end{aligned}$$

よって

$$f_{\beta, \lambda_0}(i_\alpha | i_1, \dots, i_{\alpha-1}; X_1, \dots, X_\alpha) = f_{\beta, \lambda_0}(i_\alpha | X_\alpha)$$

$$f_{\beta, \lambda_0}(X_\alpha | i_1, \dots, i_{\alpha-1}; X_1, \dots, X_{\alpha-1}) = f_{\beta, \lambda_0}(X_\alpha | X_{\alpha-1})$$

であることに注意する。明らかになっている対応が成り立つ。

$$L_c^* = \int_{\beta, \lambda_0} (i_1, \dots, i_e; R_1, \dots, R_e; t_1, \dots, t_e)$$

$$L_P(\beta) = \prod_{\alpha=1}^p \int_{\beta} (i_\alpha | \lambda_\alpha)$$

$$L_D(\beta, \lambda_0) = \prod_{\alpha=1}^p \int_{\beta, \lambda_0} (\lambda_\alpha | \lambda_{\alpha-1})$$

簡単のため  $\beta \in \mathbb{R}^1$  の場合  $n=1$  を考える。

$$U_\alpha = \partial \ln \int_{\beta} (i_\alpha | \lambda_\alpha) / \partial \beta, \quad U = \sum_{\alpha=1}^p U_\alpha = \partial \ln L_P(\beta) / \partial \beta$$

$$J(\beta) = - \partial^2 \ln L_P(\beta) / \partial \beta^2$$

$\beta^*$  :  $L_P(\beta)$  より求めらるる最尤推定量

とおく。  $L_P$  は普通の尤度関数と考え、通常最尤法が適用できるためには、この定理が、条件が与えられた分布  $\int_{\beta, \lambda_0} (i_1, \dots, i_e; R_1, \dots, R_e; t_1, \dots, t_e)$  に対し成立することを示せばよい。

定理  $n \rightarrow \infty$  のとき

$$i) \quad U / \sqrt{J(\beta^*)} \rightarrow N(0, 1) \quad \text{in law}$$

$$ii) \quad (\beta^* - \beta) \sqrt{J(\beta^*)} \rightarrow N(0, 1) \quad \text{in law}$$

を示すには、  $\int_{\beta, \lambda_0} (i_1, \dots, i_e; R_1, \dots, R_e; t_1, \dots, t_e)$  に関し、期待値  $E$  とするに  $E(U) = 0$ ,  $V(U) = \sum_{\alpha=1}^p V(U_\alpha)$  が成立すること、すなわち

$$\textcircled{1} \quad E(U_\alpha) = 0$$

$$\textcircled{2} \quad E(U_\alpha U_{\alpha'}) = 0 \quad (\alpha < \alpha')$$

が成立することを示せばよい。  $\textcircled{1}$  については、  $U_\alpha$  は  $\lambda_\alpha$  と  $\beta$  による  $n$  と  $\beta$  の条件、密度によつて定義されるから

$E(U_\alpha | \lambda_\alpha) = 0$  は明らかである。よって、 $EU_\alpha = EE(U_\alpha | \lambda_\alpha) = 0$ 。

② 同様に、 $\alpha < \alpha'$  のとき

$$E(U_\alpha U_{\alpha'}) = EE(U_\alpha U_{\alpha'} | U_\alpha = u_\alpha) = E[U_\alpha E(U_{\alpha'} | U_\alpha = u_\alpha)]$$

事象は経時的に生起するから  $\lambda_{\alpha'} \in \mathcal{F}_t$  とすると、 $\gamma$  は  $U_\alpha = u_\alpha \in \mathcal{F}_t$  とすると  $\gamma$  も含まれる。よって

$$E(U_{\alpha'} | U_\alpha = u_\alpha) = E(U_{\alpha'} | \lambda_{\alpha'}) = 0$$

よって ② は明らかである。

## 6. 部分尤度法の効率

$\beta$  に関する情報は明らかに  $L_D(\beta, \lambda_0)$  の方にも含まれている。にもかかわらず、これを考慮すると  $\beta$  の推定量は局外パラメータ  $\lambda_0$  に依存し、たまたま問題が生じるため Cox はこれを切り捨てた。しかしながら切り捨てられる情報が無視できるほど大きければ部分尤度法の適用は躊躇せざるを得ない。

$p = 1$  のとき、 $L_C^*$ ,  $L_P$  より求められる Fisher の情報量は

$$I_C^*(\beta) = E\left[-\frac{\partial^2 \ln L_C^*}{\partial \beta^2}\right] = \sum_{j=1}^q E\left\{\frac{\sum_{i \in R_j} z_i^2 e^{\beta z_i}}{\sum_{i \in R_j} e^{\beta z_i}} \int_{t_{j-1}}^{t_j} \lambda_0(t) dt\right\}$$

$$I_P(\beta) = E\left[-\frac{\partial^2 \ln L_P}{\partial \beta^2}\right] = \sum_{j=1}^q E\left\{\frac{\sum_{i \in R_j} z_i^2 e^{\beta z_i}}{\sum_{i \in R_j} e^{\beta z_i}} - \left[\frac{\sum_{i \in R_j} z_i e^{\beta z_i}}{\sum_{i \in R_j} e^{\beta z_i}}\right]^2\right\}$$

$L_C^*$  の替りに  $L_P$  を用いることより生じる情報の損失は

$$R(\beta) = \frac{I_p(\beta)}{I_c^*(\beta)}$$

で評価できる。  $R(\beta)$  を部分尤度法の効率とみなす。

特に 2 標本のとき

$$x_i = \begin{cases} 0 & : \text{標本 1} \\ 1 & : \text{標本 2} \end{cases} \quad \begin{aligned} n_0(t_j) &: R_j \text{ に属する標本 1 の回数} \\ n_1(t_j) &: R_j \text{ に属する標本 2 の回数} \end{aligned}$$

$$\lambda_0(t) = 1, \text{ 即ち指数分布のとき}$$

とすると

$$R(\beta) = \frac{\sum_{j=1}^k E\{n_1(t_j)[t_j - t_{j-1}]\}}{\sum_{j=1}^k E\left\{\frac{n_0(t_j)n_1(t_j)}{[n_0(t_j) + n_1(t_j)]e^\beta}\right\}}$$

Kalbfleisch (1974) は脱落がない場合、Efron [1977] は数学的な近似、Kay (1979) はモンテカルロ法によって  $R(\beta)$  の評価を行なったが、彼らの方法は極めて複雑、難解で筆者の力不足のため要領よく紹介できない。理論的にもまだ多くの問題点が残っており、効率の評価は今後の課題である。彼らによつて得られた結果はほぼつぎのようにまとめられる。

(i) 脱落がないとき、 $\beta = 0$  のときの評価

標本数 (N)	10	20	40	60	$\infty$
$R(0)$	0.89	0.94	0.95	0.97	1

(ii)  $R(\beta)$  は  $\beta = 0$  のとき最大値  $E$  とし  $\beta$  が増加するにつれて減少する。

(iii)  $N \rightarrow \infty$  のとき、脱落のパラメータが極端に小さければ  $\beta = 0.5$  ぐらいのとき  $R(\beta) \geq 0.80$  と考えられる。多くの場合、この値は 95% 以上になる。  $N$  が有限のときは、上の表に  $N \rightarrow \infty$  のときの値  $E$  をかけた値がほぼ  $R(\beta)$  の値に近いと見なしてよいのではないかと考えられている。

## 7. 今後の問題点

前節の2標本問題は

$$\text{標本1の分布: } F(t) = \exp\{-\int_0^t \lambda_0(x) dx\}$$

$$\text{標本2の分布: } G(t) = \exp\{-e^\beta \int_0^t \lambda_0(x) dx\}$$

であるとき、仮説  $H_0: \beta = 0$  を検定する問題にほかわらない。

$\lambda_0$  は任意だから上のように  $F(t)$  に特別の型を指定しても意味はない。したがって、問題は標本1, 標本2の母集団分布  $F, G$  とし、 $F, G$  は未知とすると、 $G = F^\psi$

( $\psi = e^\beta$ ) と表わされることから仮説:  $\psi = 1$  の検定問題に帰着する。この問題は、Lehmann<sup>(1953)</sup> の対立仮説に対する検定問題としてノンパラメトリック法では古くからよく研究されている。2標本問題に於いても同様に対処できる。従って Cox 法は脱落のあるデータに基づいて Lehmann 対立仮説を検定す



る方法とその特別の場合として提供する。他方、ノンパラメトリック法の枠内でも局所的最強の順位検定という基準に基づき脱滞のあるデータに対する順位検定が研究されてくる(例えば Johnson and Mehrota (1972), Mehrota and Johnson (1976) など)。Cox の部分尤度法は極めて広範な問題に適用可能であるが、上のように問題を制限し、ノンパラメトリック法の立場から脚光を当て、その最適性を何らかの意味で数学的に証明することは可能であるように思われる。Cox 法を数学的に明確に位置づけることは今後に残される課題がある。

### 参考文献

- Cox, D.R. (1972) : Regression models and life tables (with discussion) J.R.S.S. B34, 187-220
- Breslow, N. (1972) : 上の論文の Discussion
- Fraser, D.A.S. (1968) : The structure of inference, New York, Wiley
- Kalbfleisch, J.D. and Prentice, R.L. (1973) : Marginal likelihood based on Cox's regression and life model. Biometrika 60, 267-278
- Kaplan, E.L. and Meier, P. (1958) : Nonparametric estimation from incomplete observation J.A.S.A. 53, 457-481
- Cox, D.R. (1975) : Partial likelihood, Biometrika, 62, 269-279

Kalbfleisch, J.D. (1974) : Some efficiency calculations for survival distributions.

Biometrika 61, 31-38

Efron, B. (1977) : The efficiency of Cox's likelihood function for censored data.

J.A.S.A. 359, 557-565

Kay, R. (1979) : Some further asymptotic efficiency calculations for survival data regression models. Biometrika 66, 91-96

Lehmann, E.L. (1953) : The power of rank tests. Ann. Math. Statist. 24, 23-43

Johnson, R.A. and Mehrota, K.G. (1972) : Locally most powerful rank tests for the two-sample problem with censored data. Ann. Math. Statist. 43, 823-831

Mehrota, K.G. and Johnson, R.A. (1976) : Asymptotic sufficiency and asymptotically most powerful tests for the two sample censored situation. The Ann. of Statist. 4, 589-596

松原義弘・後藤昌司 (1979) : センサー帰による生存時間データの解析, 応用統計とコンピュータ集「医学・生物学における統計的諸問題とデータ解析」

## 付録 [コンピュータ - プログラム化へのステップ

- 重複度のある一般の場合 - ]

### (1) データの整理

死亡時刻	重複度	共変量			
$t_{ij}$	$m_{ij}$	$Z_{1ij}$	$Z_{2ij}$	...	$Z_{rij}$

の output.  $Z_{kij}$  より  $\bar{y}$  の平均  $\bar{Z}_k$  と差を引く, それを改

ために  $\sum_{k=1}^p z_{kj}$  と定めよう、 $k=1, 2, \dots, p$  .

$$(2) \quad \left. \begin{array}{l} t_{j1} \\ \text{重複) } t_{j2} \\ \vdots \\ t_{jm_j} \end{array} \right\} \begin{array}{l} \text{--- } z_{j1} \\ m_{j1} \text{ --- } z_{j2} \\ \vdots \\ \text{--- } z_{jm_j} \end{array} \quad \Delta_{ij} = \sum_{l=1}^{m_{ij}} z_{lij} \quad \text{と求める。} \quad (A1)$$

リスク集合  $R_j$  の中から相異なる  $m_{ij}$  個の要素の組を抽出し、その各々に対して、上の  $\Delta_{ij}$  を計算する。

$$R_j \left\{ \begin{array}{l} 0_{m_{j1}} \rightarrow \Delta_{j1} \\ 0_{m_{j2}} \rightarrow \Delta_{j2} \\ \vdots \\ 0_{m_{jm_j}} \rightarrow \Delta_{jm_j} \end{array} \right. \left\{ \begin{array}{l} R_j \text{ の要素の個数 } n_j \text{ とすると} \\ \text{この個数は} \\ (n_j^{ij}) \text{ 個} \end{array} \right. \quad (A2)$$

対数尤度関数は

$$\ln L_p(\beta) = \sum_{j=1}^p \beta' \Delta_{ij} - \sum_{j=1}^p \ln \left[ \sum_{k \in R_j}^* \exp \beta' \Delta_{jk} \right]$$

ここに  $\sum_{k \in R_j}^*$  は (A2) の  $(n_j^{ij})$  個にわたる和を意味する。

(3)  $\beta$  の推定、検定

$\beta = 0$  を初期値として、くり返し計算で  $\ln L_p(\beta)$  を最大にする  $\hat{\beta}$  を求める。  $\partial \ln L_p^*(\beta) / \partial \beta_k |_{\beta=\hat{\beta}}$  , Fisher の情報行列、分散共分散行列の推定値、ノンパラメトリック化された  $\hat{\beta}$  などの output。  $\beta$  に関する仮説の検定は通常の尤度比検定による。或いは AIC 基準を適用する。

(4) 生存関数の推定

(3) により 共変量の 行列  $p$  と  $\beta^* = \hat{\beta}$  を定め、それを用いて

下式に用いる。

$$\hat{\pi}_k = \frac{M_k}{\sum_{k \in R} \exp(\beta^* z_k)}$$

$$\hat{H}_0(t_{(k)}) = \prod_{j=1}^k (1 - \hat{\pi}_j)$$

$\hat{H}_0$  の分散の 1 つの近似的な推定量は次式で与えられる。

$$\begin{aligned} V[\hat{H}_0(t_{(k)})] &\doteq \left[ 1 - \frac{M_k}{\sum_{k \in R} e^{\beta^* z_k}} \right]^2 V[\hat{H}_0(t_{(k-1)})] \\ &+ [\hat{H}_0(t_{(k)})]^2 \frac{M_k}{\left( \sum_{j \in R} e^{\beta^* z_j} \right)^2}, \quad k=1, 2, \dots \end{aligned} \quad (A.1)$$

よって、初期値  $V[\hat{H}_0(t_{(0)})] = 0$  より逐次に計算することにより  $\hat{H}_0$  の標準偏差を求め、生命表形式に  $t_{(k)}$ ,  $\hat{\pi}_j$ ,  $\hat{H}_0(t_{(k)})$ ,  $E(\hat{H}_0(t_{(k)}))$ ,  $k=1, 2, \dots$  を output する。

以上 2 つの層に層別されるときは各層ごとの操作を行って  $\hat{H}_0(t_{(k)})$  を求めればよい。

[補遺] (A.1) 式の導出方法

等式の右辺に限りにおいて  $\hat{H}_0$  の分散の推定量と与えられた関数は存在しない。(A.1) の導出方法を簡単に説明しておく。本文の枠組で、重複度があるとき  $\ln L_D$  は次式で与えられる。

$$\ln L_D(\lambda_1, \dots, \lambda_k) = \sum_{j=1}^k \ln \sum_{i \in R_j} e^{\beta^* z_i} + \sum_{j=1}^k (M_j - \ln \lambda_j) - \sum_{j=1}^k \sum_{i \in R_j} \lambda_j (t_{ij} - t_{(j-1)}) e^{\beta^* z_i}$$

$$\partial \ln L_D / \partial \lambda_d = 0 \text{ かつ}$$

$$\hat{\lambda}_d = m_d / \left[ \sum_{i \in R_d} (t_{i0} - t_{i-1}) e^{\beta^* z_i} \right]$$

$$I(\lambda_d) = - E \left[ \partial^2 \ln L_D / \partial \lambda_d^2 \right] = E \left[ m_d / \lambda_d^2 \right], \quad E \left[ \partial^2 \ln L_D / \partial \lambda_d \partial \lambda_r \right] = 0$$

$d > 2$ ,  $\hat{\lambda}_1, \dots, \hat{\lambda}_d$  は漸近的に互いに独立、 $\hat{\lambda}_d$  は平均  $\lambda_d$  分散  $1/I(\lambda_d)$  の正規分布に従う。

$$\begin{aligned} \ln \hat{F}(t_{(k)}) &= \ln F(t_{(k)}) + \sum_{i=1}^k \ln \{ 1 + (t_{(i)} - t_{(i-1)}) (\hat{\lambda}_i - \lambda_i) \} \\ &\sim \ln F(t_{(k)}) + \sum_{i=1}^k (t_{(i)} - t_{(i-1)}) (\hat{\lambda}_i - \lambda_i) \end{aligned}$$

$$\therefore \hat{F}(t_{(k)}) - F(t_{(k)}) \sim F(t_{(k)}) \left\{ \sum_{i=1}^k (t_{(i)} - t_{(i-1)}) (\hat{\lambda}_i - \lambda_i) \right\}$$

$d > 2$ ,  $\hat{F}(t_{(k)})$  の分散は次式で近似される。

$$V[\hat{F}(t_{(k)})] = F^2(t_{(k)}) \sum_{i=1}^k \frac{1}{E \left\{ \left( \sum_{j \in R_i} e^{\beta^* z_j} \right)^2 \frac{1}{m_i} \right\}}$$

$$\therefore \hat{V}[\hat{F}(t_{(k)})] \equiv F^2(t_{(k)}) \sum_{i=1}^k \frac{m_i}{\left( \sum_{j \in R_i} e^{\beta^* z_j} \right)^2}$$

これは単項式の形に書ける表わすことができる。(A.1.) である。(A.1.)

で  $\beta^* = 0$  とおくと

$$\hat{V}[\hat{F}(t_{(k)})] = \left( \frac{r_k - m_k}{r_k} \right)^2 \hat{V}[\hat{F}(t_{(k-1)})] + \left[ \hat{F}(t_{(k)}) \right]^2 \frac{m_k}{r_k^2}$$

右辺第1項は Kaplan-Meier 推定量の分散の近似式であるから、  
 と一致する。第2項は Kaplan-Meier の  $r_k(r_k - m_k)$  と一致する  
 と一致する。  $r_k^2$  は  $\frac{m_k}{r_k}$  と  $\frac{m_k}{r_k}$  の積である。  $\frac{m_k}{r_k}$  は  $\frac{m_k}{r_k}$  と  $\frac{m_k}{r_k}$  の積である。  
 と一致する。  $r_k^2$  は  $\frac{m_k}{r_k}$  と  $\frac{m_k}{r_k}$  の積である。  $\frac{m_k}{r_k}$  は  $\frac{m_k}{r_k}$  と  $\frac{m_k}{r_k}$  の積である。