

複数の成長パターンを持つスギ単純同齢林における炭素固定量予測

Predicting Carbon Sequestered in an Even-Aged Sugi Forest Stand through Growth Pattern Classification

柳原 宏和・吉本 敦・二宮 嘉行

Yanagihara, H., Yoshimoto, A. & Ninomiya, Y.

キーワード: 外挿での予測評価, k -平均法, 情報量基準, 正規多変量線形回帰モデル, 成長分析, モデル選択

要約: 本研究では, 複数の成長パターンが観察される林分での炭素固定量の予測手法を提示した. その方法は以下の通りである. まず, サンプル木の材積成長データに成長曲線をあてはめパラメータの推定を行い, 推定されたパラメータの値を新たな観測値とする. この観測値に対して k -平均法に基づくクラスタリングを行いサンプル木の成長パターンを分類し, そのクラスター分割を考慮に入れた正規多変量線形回帰モデルをあてはめる. 次に, 回帰モデルのパラメータ推定値により Predictive Akaike's Information Criterion (PAIC) を計算し, PAIC に基づき最適な残存木の成長パターンの分類とクラスターの個数を決定する. 最後に, 最適な予測モデルから残存木の成長曲線のパラメータの予測値を求め, 残存木の将来的な炭素固定量の予測とその漸近 $1-\alpha$ 信頼区間を求める.

Abstract: In this paper, we presented a statistical procedure of estimating the clustered multivariate linear regression model for predicting the amount of carbon sequestered in a forest stand where there exist several growth patterns. The procedure is as follows: 1) By fitting a volume growth curve to the data of each sampled tree, parameters of the applied growth curve model are estimated. 2) By setting the estimated parameters as new observations, we

classify growth patterns of sampled trees by k -means method based on the new observations. 3) We construct a multivariate normal linear regression model with dummy variables for k -clusters. 4) Among a set of the estimated regression models with the different number of clusters, the best model is selected by minimizing the resultant predictive Akaike's information criterion (PAIC) for the remaining trees. 5) Finally, by using the best set of parameters of growth curves for the remaining trees, we predict the amount of carbon sequestered by remaining trees with its asymptotically $1-\alpha$ confidence interval.

Keywords: evaluation of prediction in interpolation, growth analysis, information criterion, k -means clustering, model selection, multivariate normal linear regression model

1. はじめに

2005年2月の京都議定書の発効に伴い森林の持つ炭素吸収機能が益々重要視され、森林が持つ炭素固定量を把握する必要性が増してきている。吸収される炭素固定量は森林の単位面積当たりの総材積量(m^3/ha)にある定数を掛けることにより概算することができるため、総材積の成長量を予測することが将来的に固定される炭素量の把握に必要不可欠になる。その際、仮に同一林分内に複数の成長パターンが存在すれば、それらの成長パターンを十分に考慮した予測を行う必要が出てくる。実際に成長パターンの識別は k -平均法(MacQueen 1967)に代表されるクラスタリング手法により可能であるが、分類対象となる林木はデータ収集のためにすでに伐採されている。すなわち将来的に固定される炭素量を予測するためには残存木の成長量を予測することが必要不可欠となる。そこで本研究は、複数の成長パターンが同時に存在する林分において、材積の成長量を予測する手法を開発し、残存木により固定される炭素量の予測を試みる。

予測の手順は以下の通りである。まず、データ収集のために伐採されたサンプル木に対し材積成長曲線をあてはめパラメータの推定を行う。次に、柳原・吉本(2005)、吉本ら(2005)に示されているように、推定されるパラメータの値を新たな観測値とし、それらに対し k -平均法に基づくクラスタリングを行い、サンプル木の成長パターンを分類する。次に新たな観測値を応答変数(response variable)とし、ダミー変数を用いることによりクラスター分割

を表現した正規多変量線形回帰モデルをあてはめる。残存木に対する最適な予測モデルの選択には、Predictive Akaike's Information Criterion (PAIC; Satoh 1997)を用いて残存木の成長パターンの分類と最適なクラスター数の決定を同時に行う。最後に、最適な予測モデルを用いて残存木の材積成長曲線のパラメータを推定し、残存木の炭素固定量の予測値とその漸近 $1-\alpha$ 信頼区間を求める。

2. サンプル木のクラスター分割の決定

今、単純同齢林内の n 個のサンプル林木から材積成長量のデータが得られたとし、 v_{il} を第 i 番目の林木の t_{il} 時点（樹齢）における観測値 ($i = 1, \dots, n, l = 1, \dots, p_l$) とする。また、 $\mathbf{v}_i = (v_{i1}, \dots, v_{ip_i})'$ を i 番目の林木に対する成長データベクトルとする。ここで、この観測値ベクトル \mathbf{v}_i に以下のような観測時点に関する非線形な成長曲線モデルを仮定する。

$$[1] \quad \mathbf{v}_i = \boldsymbol{\eta}(t_i | \boldsymbol{\theta}_i) + \boldsymbol{\varepsilon}_i \quad (i = 1, \dots, n)$$

ただし、 $\boldsymbol{\eta}(t_i | \boldsymbol{\theta}_i)$ は各林木に対し q 個の要素からなる未知パラメータベクトル $\boldsymbol{\theta}_i = (\theta_{i1}, \dots, \theta_{iq})'$ を持ち、時点ベクトル $t_i = (t_{i1}, \dots, t_{ip_i})'$ の関数で定義される $p_i \times 1$ 平均値ベクトルである。また、 $\boldsymbol{\varepsilon}_i = (\varepsilon_{i1}, \dots, \varepsilon_{ip_i})'$ は、それぞれ独立に平均 $E[\boldsymbol{\varepsilon}_i] = \mathbf{0}_{p_i}$ 、分散共分散行列 $\text{Cov}[\boldsymbol{\varepsilon}_i] = \sigma_i^2 \mathbf{I}_{p_i}$ を持つ分布に従う誤差ベクトルとする。なお、 $\boldsymbol{\eta}(t_i | \boldsymbol{\theta}_i)$ は仮定する既知関数 $f(t_{ij} | \boldsymbol{\theta}_i)$ により以下のように定義するものとする。

$$[2] \quad \boldsymbol{\eta}(t_i | \boldsymbol{\theta}_i) = \begin{pmatrix} f(t_{i1} | \boldsymbol{\theta}_i) \\ \vdots \\ f(t_{ip_i} | \boldsymbol{\theta}_i) \end{pmatrix}$$

仮に、柳原・吉本(2005)で用いたリチャーズの成長関数(Richards 1958)を用いれば、

$$[3] \quad f(t_{il} | \boldsymbol{\theta}_i) = e^{\theta_{i1}} \left\{ 1 - \exp(-e^{\theta_{i2}} t_{il}) \right\}^{\exp(\theta_{i3})}$$

となり、未知パラメータの個数は $q = 3$ となる。次に[1]式の推定量 $\hat{\boldsymbol{\theta}}_i$ を推定し、柳原・吉本(2005)や吉本ら(2005)と同様に、得られる推定値 $\hat{\theta}_{i1}, \dots, \hat{\theta}_{iq}$

を新たな観測値 $y_i = (\hat{\theta}_{i1}, \dots, \hat{\theta}_{iq})'$ とみなし, $\mathbf{Y} = (y_1, \dots, y_n)'$ に対して成長パターンの分類を行う.

成長パターン分類については以下の通りである. 林分内には k 個の成長パターン (クラスター) が存在し, そのクラスター分割を $G = \{C_1, \dots, C_k\}$ とする. ただし, 第 g 番目のクラスターに属する固体の数は n_g とする. なお $n_1 + \dots + n_k = n$ である. このとき, 各クラスター C_g ($g = 1, \dots, k$) の重心 \bar{y}_g と共分散行列 S_g は,

$$[4] \quad \bar{y}_g = \frac{1}{n_g} \sum_{i \in C_g} y_i, \quad S_g = \frac{1}{n_g} \sum_{i \in C_g} (y_i - \bar{y}_g)(y_i - \bar{y}_g)'$$

である. MacQueen(1967)により提案された k -平均法は, 以下のように分散分析における群内平方和と類似したクラスター内平方和積和行列 $\mathbf{W}(G) = n_1 S_1 + \dots + n_k S_k$ の変化量の増減により個体 i が属するクラスターを決定する. 第 i 番目の個体がクラスター C_g に属しているようなクラスター分割 $G = \{C_1, \dots, C_k\}$ において, この個体を他のクラスター C_h ($h \neq g$) に移動させるとする. このとき, 新しいクラスター分割 $G^* = \{C_1, \dots, C_h^*, \dots, C_g^*, \dots, C_k\}$ におけるクラスター内平方和積和行列 $\mathbf{W}(G^*)$ を計算し, 個体 i をクラスター C_h に移動させる前の分割 G に基づく $\mathbf{W}(G)$ に比べ $\mathbf{W}(G^*)$ が小さくなっていけば分割を更新し, そうでなければ更新しない. このようにクラスター分割の更新を逐次判定し, 分割が収束するまで続ける.

上記のようにクラスター分割を更新する前のクラスター内平方和積和行列 $\mathbf{W}(G)$ と分割を更新した後でのクラスター内平方和積和行列 $\mathbf{W}(G^*)$ の差の増減により分割を更新するか否かを決定する手法が k -平均法であるが, 実際 $\mathbf{W}(G)$ と $\mathbf{W}(G^*)$ は行列であるため, 大きさの測り方には様々な基準がある. 一般的な手法は行列のトレースにより大きさを測り,

$$[5] \quad \text{tr}(\mathbf{W}(G)) > \text{tr}(\mathbf{W}(G^*))$$

であればクラスター分割を更新するというものである. しかしながら, このトレースによる更新判定は単純に個体 i とクラスター重心とのユークリッド距離の大小で行われ, 同一個体内の変数同士が強い相関を持つ場合, この判定基準では相関を考慮したクラスター分割が困難となる. すなわちトレース

のみで大きさを比較する場合、単純にユークリッド距離に近いもの同士を同一のクラスターと識別するからである。そこで本論文では強い相関を考慮できるように、行列式により大きさを判断する手法を用いる。すなわち、

$$[6] \quad |W(G)| > |W(G^*)|$$

であればクラスター分割を更新するという手法である。

クラスター分割を更新するか否かを判断するためには、逐次更新される新しい分割でのクラスター内平方和積和行列 $W(G^*)$ を求める必要がある。しかしながら、柳原・吉本(2005)に示されているように、それぞれの分割におけるクラスター内平方和積和行列 $W(G)$ と $W(G^*)$ の行列式の比、

$$[7] \quad \frac{|W(G^*)|}{|W(G)|} \text{は、} \quad a_g = \sqrt{\frac{n_g}{n_g - 1}} W(G)^{-1/2} (y_i - \bar{y}_g) \quad a_h = \sqrt{\frac{n_h}{n_h + 1}} W(G)^{-1/2} (y_i - \bar{y}_h)$$

とおくと、

$$[8] \quad \frac{|W(G^*)|}{|W(G)|} = (1 + a_h' a_h)(1 - a_g' a_g) + (a_g' a_h)^2$$

と表すことができるため、基となる分割から得られる情報のみにより分割の更新判別が可能となる。すなわち不等式 $|W(G)| > |W(G^*)|$ と $(1 + a_h' a_h)(1 - a_g' a_g) + (a_g' a_h)^2 < 1$ は同値となるからである。なお、 \bar{y}_g と \bar{y}_h は [4] 式で与えられる、分割更新前のクラスター C_g と C_h の重心である。

上記の行列式に基づく k -平均法によるクラスタリングの手順は、以下のようになる。

行列式を用いた k -平均法

Step 1: クラスターの個数 k をあらかじめ与え、乱数により各個体の分類し、初期分割を決定する。

Step 2: 第 i 番目の個体が属するクラスターを C_g とし、クラスター C_g から C_h ($h \neq g$) に移動させる状況を考える。このとき、

$$[9] \quad (1 + a_h' a_h)(1 - a_g' a_g) + (a_g' a_h)^2 < 1$$

であれば個体 i の属するクラスターを C_g から C_h に更新する。ただし a_g と a_h は [7] 式で与えられるベクトルである。

Step 3: Step 2 の作業を全ての個体に対し、すべてのクラスターについて取

束するまで行い、最終的にそれぞれが属するクラスターを決定する。

Step 4: 初期分割を変え十分にStep 1 - Step 3を繰り返し、それぞれ得られた分割でのクラスター内平方和積和行列の行列式 $|\mathbf{W}(G)|$ が最小となるクラスター分割を最終的に最適な分割とする。

k -平均法はクラスター内平方和積和行列の行列式が最小になるようにクラスター分割 G を更新する手法であり、必ず極小値に収束するが、必ずしも最小値である保証がない。そこで、初期分割を変え十分繰り返すStep 4が必要不可欠となる。

3. 伐採木のクラスター分割と最適なクラスター数の決定法

推定量 \mathbf{Y} について k -平均法を適用し k 個のクラスター分割 $G = \{C_1, \dots, C_k\}$ が得られているとする。クラスター毎に成長パターンが異なるため、クラスター毎に異なる線形モデルをあてはめる必要がある。すなわち、個体 i がクラスター C_g に属するとき、 \mathbf{y}_i に以下のような $r \times 1$ 説明変数ベクトル $\mathbf{x}_i = (x_{i1}, \dots, x_{ir})'$ を持つ正規多変量線形回帰モデルをあてはめることになる。

$$[10] \quad M_k : \mathbf{y}_i \sim i.d. N_q(\Xi_g' \mathbf{x}_i, \Sigma) \quad (i = 1, \dots, n; g = 1, \dots, k)$$

ここで Ξ_g と Σ はそれぞれ $r \times q$ と $q \times q$ の平均構造と分散共分散行列を表す未知パラメータ行列であり、 Ξ_g の成分は以下の通りである。

$$[11] \quad \Xi_g = \begin{pmatrix} \xi_{11}^g & \cdots & \xi_{1q}^g \\ \vdots & \ddots & \vdots \\ \xi_{r1}^g & \cdots & \xi_{rq}^g \end{pmatrix}$$

次に $\Xi' = (\Xi_1', \dots, \Xi_k')$ とし、 Ξ を用いてクラスター毎に異なるモデルに対しダミー変数を用いて一つのモデルで表す。今、個体 i がクラスター C_g に属せば1、属さないのであれば0であるようなダミー変数 d_{ig} を考え、そのベクトルを $\mathbf{d}_i = (d_{i1}, \dots, d_{ik})'$ とする。このとき、個体が g 番目のクラスターに属するため、 \mathbf{d}_i は g 番目の成分が1で残りが0となる $k \times 1$ ベクトルである。このとき、ダミー変数を加えて新たに $kr \times 1$ 説明変数ベクトル

$\mathbf{x}_{d,i} = \mathbf{d}_i \otimes \mathbf{x}_i$ を考える. ただし “ \otimes ” はクロネッカー積を表す. なお, $n \times m$ 行列 $\mathbf{A} = [a_{ij}]$ と行列 \mathbf{B} のクロネッカー積は,

$$[12] \quad \mathbf{A} \otimes \mathbf{B} = \begin{pmatrix} a_{11}\mathbf{B} & \cdots & a_{1m}\mathbf{B} \\ \vdots & \ddots & \vdots \\ a_{n1}\mathbf{B} & \cdots & a_{nm}\mathbf{B} \end{pmatrix}$$

である(参照 Magnus and Neudecker 1999, p.27-28). 従って $\mathbf{x}_{d,i}$ は,

$$[13] \quad \mathbf{x}_{d,i} = \mathbf{d}_i \otimes \mathbf{x}_i = \begin{pmatrix} \mathbf{0}_{(g-1)r} \\ \mathbf{x}_i \\ \mathbf{0}_{(k-g)r} \end{pmatrix}$$

となる. この説明変数ベクトル $\mathbf{x}_{d,i}$ を用いると,

$$[14] \quad \Xi' \mathbf{x}_{d,i} = \begin{pmatrix} \Xi_1' & \cdots & \Xi_g' & \cdots & \Xi_k' \end{pmatrix} \begin{pmatrix} \mathbf{0}_r \\ \vdots \\ \mathbf{x}_i \\ \vdots \\ \mathbf{0}_r \end{pmatrix} = \Xi_g' \mathbf{x}_i$$

となることから, [10] 式は

$$[15] \quad M_k : \mathbf{y}_i \sim i.d. N_q(\Xi' \mathbf{x}_{d,i}, \Sigma) \quad (i = 1, \dots, n)$$

と表すことができる. このとき, モデル M_k における \mathbf{y}_i の確率密度関数は,

$$[16] \quad L(\mathbf{y}_i | \mathbf{x}_{d,i}, \Xi, \Sigma) = \left(\frac{1}{2\pi}\right)^{q/2} |\Sigma|^{-1/2} \exp\left\{(\mathbf{y}_i - \Xi' \mathbf{x}_{d,i})' \Sigma^{-1} (\mathbf{y}_i - \Xi' \mathbf{x}_{d,i})\right\}$$

である(参照 塩谷 1990, p.7-12). 今, $n \times kr$ 行列 \mathbf{X}_d を $\mathbf{X}_d = (\mathbf{x}_{d,1}, \dots, \mathbf{x}_{d,n})'$ と置くと, モデル M_k の下での Ξ と Σ の最尤推定量は,

$$[17] \quad \hat{\Xi}_Y = (\mathbf{X}_d' \mathbf{X}_d)^{-1} \mathbf{X}_d' \mathbf{Y}, \quad \hat{\Sigma}_Y = \frac{1}{n} \mathbf{Y}' \left[\mathbf{I}_n - \mathbf{X}_d (\mathbf{X}_d' \mathbf{X}_d)^{-1} \mathbf{X}_d' \right] \mathbf{Y}$$

となる.

実際に林木の成長データを用いる際, 上記で求められるモデル M_k は伐採されたサンプル木, つまり実測値に基づくモデルであり, 残存木に関するモデルではない. すなわち残存木のクラスター分割を評価するためには, M_k のモデル評価ではなく残存木に関するモデル評価を行わなければならない.

そのためには、残存木の成長曲線のパラメータに関する予測モデルを次のように考える必要がある。

今、残存木の本数を m として、 ζ_j と δ_j ($j=1, \dots, m$) を残存木の材積成長曲線のパラメータベクトルとクラスター分割を表現する観測されていないダミー変数ベクトルとする。また観測される残存木の説明変数ベクトルを w_j とし、 $w_{\delta,j} = \delta_j \otimes w_j$ とする。このとき、残存木もサンプル木から得られる実測値モデル M_k と同様のプロセスに従うと仮定し、残存木に関する予測モデルを

$$[18] \quad M_k : \zeta_j \sim i.d. N_q(\Xi' w_{\delta,j}, \Sigma) \quad (j=1, \dots, m)$$

とする。なお、予測モデル M_k において、 Ξ と Σ は実測値モデル M_k と共通である。

最適な実測値モデル M_k の選択には予測カルバックライブラーの距離 (Kullback and Leibler 1951) に基づくリスクを用いることが一般的である。そのリスクは、 \mathbf{Y} と独立で同一な分布に従う確率変数 $\mathbf{Y}_F = (y_{F,1}, \dots, y_{F,m})'$ を用いて、

$$[19] \quad R_k(\mathbf{X}_d | \mathbf{X}_d) = -2 \sum_{i=1}^n E_Y^* E_{Y_F}^* \left[\log L(y_{F,i} | \mathbf{x}_{d,i}, \hat{\Xi}_Y, \hat{\Sigma}_Y) \right]$$

と定義される。ただし E^* は真のモデルの下での期待値である。実測値モデルにおける予測リスクに関しては、吉本ら(2005)に詳しい。上記のリスクでは、 \mathbf{X}_d という説明変数が与えられたとき、その \mathbf{X}_d の下で予測を行ったという状況でモデルを評価する。すなわちモデルの内挿での予測に基づく評価である。しかしながら、ここで評価するモデルは M_k であり、この場合、 $\mathbf{Z} = (\zeta_1, \dots, \zeta_m)'$ は観測されないため、 Ξ と Σ の推定は \mathbf{Y} から行わなくてはならない。その結果、予測に基づくリスクは、

$$[20] \quad R_k(\mathbf{W}_\delta | \mathbf{X}_d) = -2 \sum_{j=1}^m E_Y^* E_Z^* \left[\log L(\zeta_j | w_{\delta,j}, \hat{\Xi}_Y, \hat{\Sigma}_Y) \right]$$

となる。ただし、 $\mathbf{W}_\delta = (w_{\delta,1}, \dots, w_{\delta,m})'$ である。このリスクでは、 \mathbf{X}_d と異なる \mathbf{W}_δ の下で予測を行った場合のモデルを評価することになり、モデルの外挿での予測に基づくモデル評価である。仮に $\mathbf{W}_\delta = \mathbf{X}_d$ となれば、[19] 式と

[20] 式は一致する．以下ではリスク $R_k(\mathbf{W}_\delta | \mathbf{X}_d)$ の最小化により最適な予測モデル \mathcal{M}_k の選択を行う．

上記の [20] 式のリスクの推定量は Satoh(1997)により提案されたPredictive Akaike's Information Criterion (PAIC)である．予測モデル \mathcal{M}_k に関する PAIC は以下のように定義される．

$$[21] \quad \text{PAIC}_k = m \log |\hat{\Sigma}_Y| + m q \log 2\pi + \frac{n(m + \phi_k)q}{n - kr - q - 1}$$

ただし,

$$[22] \quad \phi_k = \text{tr} \left\{ \mathbf{W}_\delta' \mathbf{W}_\delta (\mathbf{X}_d' \mathbf{X}_d)^{-1} \right\} = \sum_{j=1}^m (\delta_j' \otimes \mathbf{w}_j') (\mathbf{X}_d' \mathbf{X}_d)^{-1} (\delta_j \otimes \mathbf{w}_j)$$

ここで、クラスター数 k を固定して考えると、 $\hat{\Sigma}_Y$ と $\hat{\Sigma}_Y$ は一定となり、PAIC の最小化は ϕ_k の最小化に等しくなる．このとき ϕ_k の値は残存木のクラスター分割 $\delta_1, \dots, \delta_m$ のみによって決定されるため、 ϕ_k を最小にする $\delta_1, \dots, \delta_m$ はそのまま PAIC を最小にするクラスター分割となる．よって、 k を固定した下での残存木の最適なクラスター分割は、 ϕ_k を最小にする $\delta_1, \dots, \delta_m$ になる．個々の $(\delta_j' \otimes \mathbf{w}_j') (\mathbf{X}_d' \mathbf{X}_d)^{-1} (\delta_j \otimes \mathbf{w}_j)$ は正定値行列の二次形式であり、その値は常に正になることから、 ϕ_k を最小にするためには $(\delta_j' \otimes \mathbf{w}_j') (\mathbf{X}_d' \mathbf{X}_d)^{-1} (\delta_j \otimes \mathbf{w}_j)$ を最も小さくするように δ_j を決定すればよいことになる．すなわち、最適なクラスター分割は、 g 番目の成分が 1 で残りが 0 であるような $k \times 1$ ベクトル e_g を用いると、

$$[23] \quad \hat{\delta}_j = \arg \min_{e=e_1, \dots, e_k} (e' \otimes \mathbf{w}_j') (\mathbf{X}_d' \mathbf{X}_d)^{-1} (e \otimes \mathbf{w}_j) \quad (j=1, \dots, m)$$

となる．更に、残存木に対するクラスターの個数を変えてそれぞれのクラスターの個数でのPAICを比較することにより最適なクラスター数も決定する．つまり、それぞれの k の下でPAICが最小となるクラスター分割 $\hat{\delta}_1, \dots, \hat{\delta}_m$ を決定し、その分割の下でのPAICを求め、PAICを最小とする k を最適なクラスター数とみなす．以下に残存木に対するクラスター分割と最適なクラスター数を決定する手順を示す．

残存木のクラスター分割と最適なクラスターの個数の決定法

Step 1: 最大のクラスター数 K を決め、 $k=1, \dots, K$ と繰り返す．

Step 2: 伐採された個々のサンプル木の材積成長データに成長曲線をあては

め、その推定された係数を新しい観測値 \mathbf{Y} とする。

Step 3: クラスターの個数 k に対し、 \mathbf{Y} にクラスター内平方和積和行列の行列式の最小化による k -平均法を適用し、サンプル木の最適なクラスター分割 d_1, \dots, d_m を決定する。

Step 4: 説明変数行列を $\mathbf{X}_d = (d_1 \otimes \mathbf{x}_1, \dots, d_m \otimes \mathbf{x}_m)'$ とし、 $\hat{\Xi}_Y$ と $\hat{\Sigma}_Y$ を求め、 ϕ_k の最小化により残存木のクラスター分割 $\hat{\delta}_1, \dots, \hat{\delta}_m$ を決定する。これらを用いて PAIC を計算する。

Step 5: Step 3 と 4 を繰り返し、以下のように PAIC を最小にする k を最適なクラスター数 k_{opt} とする。

$$[24] \quad k_{\text{opt}} = \arg \min_{k=1, \dots, K} \text{PAIC}_k$$

4. 炭素固定量の予測

ここでは、 n 本の林木の材積成長データが得られている面積 S (ha) を持つ試験林に、 m 本の林木が伐採されずに残っていると、この残存木の炭素固定量を予測する。スギ林分の炭素固定量は以下のような単位あたりの総材積により求める(参照 松本 2001)。

$$\text{炭素固定量 (Ct/ha)} = \text{総幹材積 (m}^3\text{)} \times 0.38 (\text{比重})$$

$$[25] \quad \times 0.44 (\text{スギ炭素割合}) \times \frac{100}{61} \times \frac{1}{\text{林分面積 (ha)}}$$

すなわち、林分内の残存木の総材積量を予測できれば、炭素固定量も同時に予測することができる。

まず、残存木での総材積はそれぞれの林木の材積成長曲線の総和によって予測する。このとき、残存木の材積成長曲線のパラメータを前章と同様に ζ_j ($j=1, \dots, m$) と置くと、時点 t における炭素固定量の予測値は、

$$[26] \quad G_{\text{cs}}(t) = \frac{16.72}{61S} \sum_{j=1}^m f(t | \zeta_j) \quad (\text{Ct/ha})$$

となる。しかしながら、 ζ_j は観測されない変数であり、その値を実際に用いることはできないため、予測モデルにより ζ_j を予測し、その予測値を用いて [26] 式の推定量を求める。今、 $\hat{\Xi}_Y$ と $\hat{\Sigma}_Y$ を PAIC の最小化により選ばれた最適なモデルの Ξ と Σ の推定量、 $\hat{\delta}_1, \dots, \hat{\delta}_m$ を残存木の最適なモデルでの

最適なクラスター分割とする。このとき $w_{\delta_j} = \hat{\delta}_j \otimes w_j$ とすると、 ζ_j の予測値は $\hat{\zeta}_j = \hat{\Xi}_Y' w_{\delta_j}$ となる。 $G_{cs}(t)$ の推定量は、[26] 式で観測されない ζ_j を予測値 $\hat{\zeta}_j$ で置き換えることにより求めることができ、その推定量は以下のようになる。

$$[27] \quad \hat{G}_{cs}(t) = \frac{16.72}{61S} \sum_{j=1}^m f(t | \hat{\Xi}_Y' w_{\delta_j}) \quad (\text{Ct/ha})$$

この推定量 $\hat{G}_{cs}(t)$ の値は実測値 \mathbf{Y} によって変動するため、推定値の変動を考慮した推定も行うことが望ましい。その場合、信頼係数 $1-\alpha$ の区間推定を行えばよいが、成長曲線モデルが非線形関数であるため、信頼係数が正確に $1-\alpha$ となる信頼区間を構成することはできない。そのため、信頼係数は漸近的なものによって代用することにする。今、 z_α を標準正規分布の上側 $100 \times \alpha \%$ 点とし、偏微分に関する $q \times 1$ ベクトル $h_j(t | \Xi_0)$ と $krq \times 1$ ベクトル $g_j(t | \Xi_0)$ を以下のように定義する。

$$[28] \quad \begin{aligned} h_j(t | \Xi_0) &= \left. \frac{\partial}{\partial \zeta} f(t | \zeta) \right|_{\zeta = \Xi_0' w_{\delta_j}} \\ g_j(t | \Xi_0) &= \left. \frac{\partial}{\partial \text{vec}(\Xi)} f(t | \Xi' w_{\delta_j}) \right|_{\Xi = \Xi_0} = (\mathbf{I}_q \otimes w_{\delta_j}) h_j(t | \Xi_0) \end{aligned}$$

このとき、

$$[29] \quad \bar{g}(t | \Xi_0) = \frac{1}{m} \sum_{j=1}^m g_j(t | \Xi_0) \quad \bar{H}(t | \Xi_0) = \frac{1}{m} \sum_{j=1}^m h_j(t | \Xi_0) h_j(t | \Xi_0)'$$

とし、

$$[30] \quad \psi^2(t | \Xi_0, \Sigma_0) = \text{tr} \left\{ \bar{H}(t | \Xi_0) \Sigma_0 \right\} + m \bar{g}(t | \Xi_0)' \left[\Sigma_0 \otimes (\mathbf{X}_d' \mathbf{X}_d)^{-1} \right] \bar{g}(t | \Xi_0)$$

とする。ただし、 $\text{vec}(\mathbf{A})$ は行列 \mathbf{A} の列ベクトルを縦に並べてできるベクトルを表し、 $n \times m$ 行列 $\mathbf{A} = [a_{ij}]$ であれば、

$$[31] \quad \text{vec}(\mathbf{A}) = (a_{11}, a_{21}, \dots, a_{n1}, \dots, a_{1m}, a_{2m}, \dots, a_{nm})'$$

となる(参照 Magnus and Neudecker 1999, p. 30-31)。このとき $G_{cs}(t)$ の漸近 $1-\alpha$ 信頼区間を以下のように定義する。

$$[32] \quad \hat{G}_{cs}^-(t) \leq G_{cs}(t) \leq \hat{G}_{cs}^+(t)$$

ただし、

$$\begin{aligned}
 \hat{G}_{cs}^-(t) &= \hat{G}_{cs}(t) - \left(\frac{16.77}{61S}\right) \sqrt{m} z_{\alpha/2} \psi(t | \hat{\Xi}_Y, \hat{\Sigma}_Y) \quad (Ct/ha) \\
 \hat{G}_{cs}^+(t) &= \hat{G}_{cs}(t) + \left(\frac{16.77}{61S}\right) \sqrt{m} z_{\alpha/2} \psi(t | \hat{\Xi}_Y, \hat{\Sigma}_Y) \quad (Ct/ha)
 \end{aligned}
 \tag{33}$$

であり, $\psi(t | \Xi_0, \Sigma_0) = \sqrt{\psi^2(t | \Xi_0, \Sigma_0)}$ である. このとき以下のような定理が成り立つ. なお, 定理の証明については Appendix に記す.

定理: $\sqrt{n/m} = O(1)$ を仮定する. このとき,

$$P(\hat{G}_{cs}^-(t) \leq \mathcal{G}_{cs}(t) \leq \hat{G}_{cs}^+(t)) \rightarrow 1 - \alpha \quad (n \rightarrow \infty)$$

が成り立つ.

5. 実データへの適用

本章では, 前章で導入した手法を用いて, 炭素固定量の予測を行った. 今回解析に使用したデータは, 柳原・吉本(2005)で用いた福岡県八女郡星野村における23年生の無間伐林より抽出した30(= n)本のサンプル木から得た成長データである. 残存木の本数は106(= m)本であった. 成長曲線としては, [3]式でのリチャーズの成長関数を用い, その推定値を $y_i = (\hat{\theta}_{i1}, \hat{\theta}_{i2}, \hat{\theta}_{i3})'$ ($i = 1, \dots, 30$)として解析に用いた. 試験林の形や立木位置, また材積成長データと成長曲線のパラメータの推定値等は, 柳原・吉本(2005)に詳しい. 本解析において使用した説明変数は胸高直径DBHである. 従って説明変数ベクトルはそれぞれ $x_i = (1, \text{DBH}_i)'$ ($i = 1, \dots, 30$), $w_j = (1, \text{DBH}_j)'$ ($j = 1, \dots, 106$)となり, $r = 2$ となる. 図1に試験林でのDBH

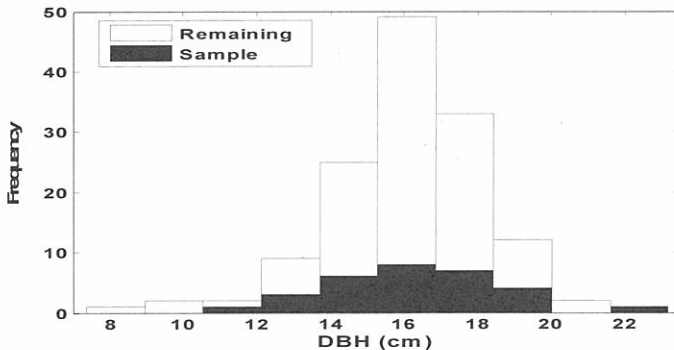


図1. サンプル木と残存木のDBHのヒストグラム

のヒストグラムを記す．それぞれ黒く塗りつぶれた棒がサンプル木のDBH，白抜き棒が残存木のDBHを表している．この図から分布の形状は異なるものの，レンジに関してはサンプル木と伐採木でさほど変わらないことがわかる．

まずサンプル木に対して，最大クラスター数 $K=5$ として行列式に基づく k -平均法を用いてクラスタリングを行った．図2.1, 2.2, 2.3, 2.4は $k=2, \dots, 5$ でのクラスター分割の結果である．図中の n_g ($g=1, \dots, k$) はクラスター C_g に属していると判断されたサンプル木の本数を表している．なお $n_1 + \dots + n_k = n$ である．これらの図から，どの分割においても共通の 8 本の林木がクラスター 1 に分類されていることがわかる．そのため，これらの 8 本の林木は他の 22 本とは明らかに異なる成長パターンを持つと推測される．その他 22 本のクラスター分割は k の個数によって異なるが，クラスターの個数を増やせば，分割されたクラスターに属する標本の数はほぼ同じになることがわかる．

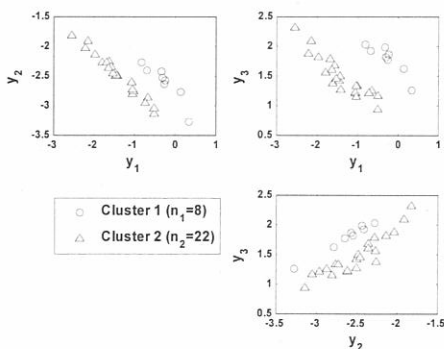


図2.1. サンプル木でのクラスタリング結果 ($k=2$)

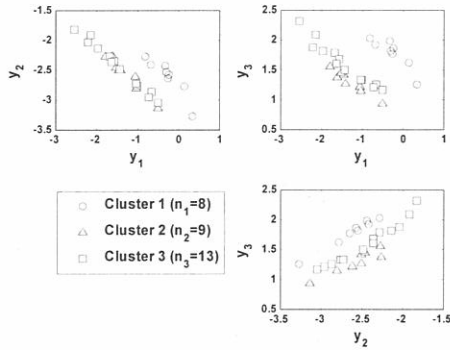


図2.2. サンプル木でのクラスタリング結果 ($k=3$)

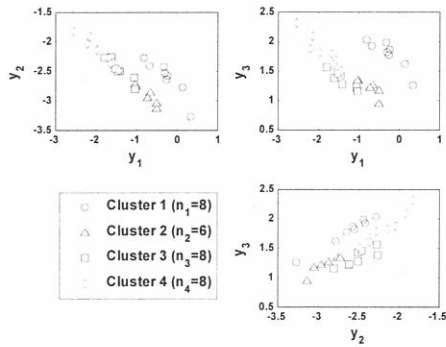


図2.3. サンプル木でのクラスタリング結果 ($k=4$)

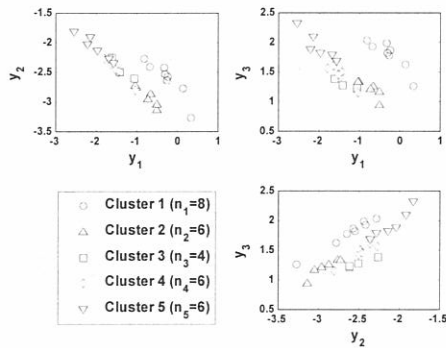


図2.4. サンプル木でのクラスタリング結果 ($k=5$)

次に、サンプル木のクラスタリング結果に基づき正規多変量線形モデルをあてはめ、予測モデル $\{M_1, \dots, M_k\}$ に対するPAICを比較することにより最適なクラスター数と予測モデルを決定した。それぞれのモデルに対するPAICは図3のようになり、3個のクラスターを持つモデルが最適なモデルとして選ばれた。大まかに見ると、観測値 Y のクラスターは2個のように観察されるが、 $y_1 - y_2$ や $y_1 - y_3$ での散布図では不明であるものの、 $y_2 - y_3$ での散布図を見ると、3個目のクラスターがあるように観察できる。PAICの値から、この3個目のクラスターもモデル化した方が良いと判断したことになる。

最適なモデルからの残存木の材積成長曲線のパラメータの予測結果を図4に示す。図中の m_g ($g = 1, \dots, k$) はクラスター C_g に属していると判断された残存木の本数を表している。なお $m_1 + \dots + m_k = m$ である。サンプル木の場合と比べると、3番目のクラスターに属するデータが長く伸びている。これは、サンプル木では観測されなかった非常に小さいDBHが残存木にあるからである。それらを除けばサンプル木でのクラスターとはほぼ同様のクラスター分割が得られる。

残存木のクラスター分類とDBHの大きさの関係を図5に示す。この図から、大まかにDBHが比較的大きい林木は1番目のクラスターに、平均的な林木は2番目のクラスターに、そして比較的小さい林木は3番目のクラスターに分類されていることがわかる。

最後に、ここで得られた予測モデルを用いて計算された炭素固定量の予測結果を図6に示す。図6は1から5個のクラスターを持つモデルにおいて計算された $\hat{G}_{cs}(t)$ に対し5年おきに求めたグラフである。また図7は最適なモデル $k=3$ での残存木の炭素固定量の漸近0.95信頼区間を表している。図6から、クラスターの数が増加するに伴い、炭素固定量の予測値が増えているのがわかる。その結果、もともと成長パターンが異なる林木に対して、パターンの違いを考慮しない場合、推定値の当てはまりが悪くなり、予測値を過小評価する傾向があると予想できる。逆にクラスターの個数を多めに設定すれ

ば、クラスター分類に関して誤判別を起こす危険性が高まり、予測値が増えたとしてもその値の信頼性は低くなる可能性がある。

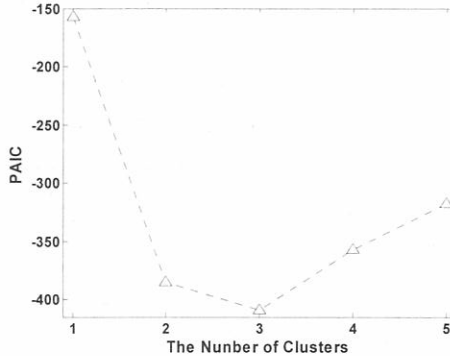


図3. クラスターの個数とPAIC

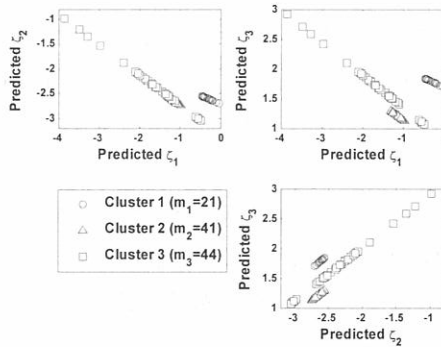


図4. 最適な予測モデルでの残存木の成長パラメータの予測値

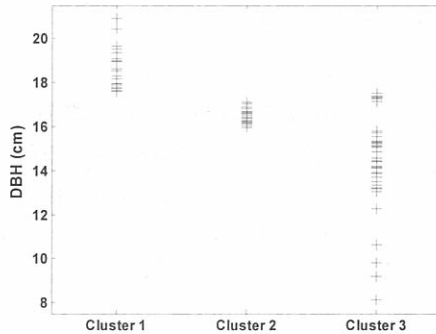


図5. 最適モデルでのクラスター分類とDBHの関係

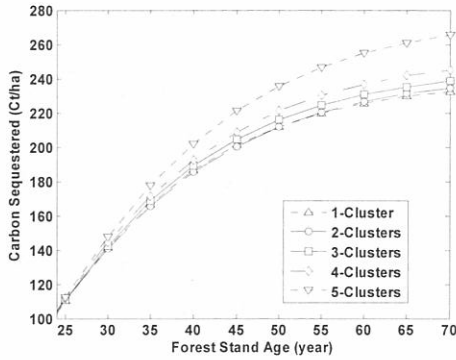


図6. それぞれのクラスターでの炭素固定量の予測量

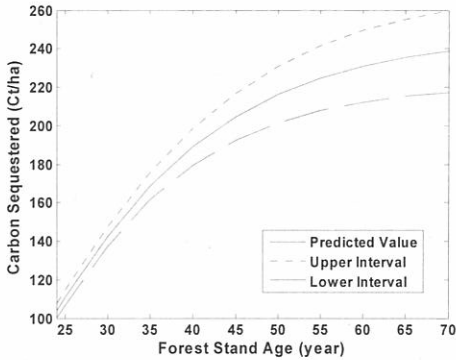


図7. 最適なモデルでの炭素固定量の予測量

6. おわりに

本論文では、同一林分の内に複数の成長パターンが存在するとき、サンプル木から残存木の成長パターンを分類し、その分割に基づいて炭素固定量を予測する手法を提示した。ここで用いた方法では、実際に残存木の成長データが得られないため、外挿での予測に基づくリスクの推定量であるPAICの最小化により、残存木のクラスター分割と最適なクラスター数、及び最適な予測モデルを決定した。ここで提示された手法と実データを用いて、残存木の成長曲線のパラメータを予測し、その値を用いて将来的に固定される炭素量の予測及びその漸近0.95信頼区間の推定を行った。

近年、一般化非線形混合効果モデル(Generalized Non-linear Mixed-effects Model; Vonesh and Carter 1992)等の多変量モデルを用いた成長データの開発が行われている。また、その種が多変量モデルを用いた炭素固定量の予測も行われている(Yanagihara and Yoshimoto 2005)。しかしながら、それらのモデルで残存木の成長パターンの分類を行える有効な手法は開発されていない。また混合効果モデルはモデルの特定のために莫大な計算量を必要とするため、候補のモデルが多い場合、計算量的にもその使用は困難である。そのような問題点に対し、今回提示した手法は簡便でかつ有効な手法であると考えられる。

実データを用いた分析では、説明変数としてDBHのみを用いたが、成長に影響を与える因子は他にも考えられる。吉本ら(2005)は、線形回帰モデルにおいて、周辺のDBHの平均、特に、半径5 m以内の周りの林木のDBHの重み付き平均もDBHと同様に最適な説明変数として選択している。仮にそれらの変数を用いて内挿での予測精度が上がれば、外挿での予測精度もあがる。そのため、より精度の高い予測ができるような説明変数の組み合わせを探索し、炭素固定量の予測を行うことも今後の課題である。

謝辞

本研究は、文部科学省科学研究費（基盤研究（B）(2):No.15330048）及び地球環境研究総合推進費S-4を受けて行われたものである。

引用文献

- Kullback, S. and Leibler, R. A. 1951. On information and sufficiency, *Annal of Mathematical Statistics* 22:79-86
- MacQueen, J. B. 1967. Some methods for classification and analysis of multivariate observations, pp.281-298, Proceedings of the Fifth Berkeley Symposium on Mathematical Statistics and Probability (Neyman, J. ed.), 1, Berkeley,
- Magnus, J. R. and Neudecker, H. 1999. Matrix differential calculus with applications in statistics and econometrics (Revised ed.), John Wiley & Sons, New York

- 松本光朗. 2001. 日本の森林による炭素蓄積量と炭素吸収量, *森林科学* 33:30-36
- Richards, F. J. 1958. A flexible growth function to empirical use, *Journal of Experimental Botany* 10:290-300
- Satoh, K. 1997. AIC-type model selection criterion for multivariate linear regression with a future experiment, *Journal of the Japan Statistical Society* 27:135-140
- Seber, G. A. F. and Wild, C. J. 1989. *Nonlinear Regression*, John Wiley & Sons, New York
- 塩谷實. 1990. *多変量解析概論*, 朝倉書店, 東京
- Vonesh, E. F. and Carter, R. L. 1992. Mixed-effects nonlinear regression for unbalanced repeated measures, *Biometrics* 48:1-17
- 柳原宏和・吉本敦. 2005. 単純同齢林における林木成長パターンのクラスタリング, pp.49-69, *森林資源管理と数理モデル Vol. 4* (近藤洋史・吉本敦・松村直人 編集), 森林計画学会出版局, 東京
- Yanagihara, H. and Yoshimoto, A. 2005. Statistical procedure for assessing the amount of carbon sequestered by sugi (*Cryptomeria japonica*) plantation, pp.125-140, *Multipurpose Inventory for the Aged Artificial Forest* (Nobori, Y., Takahashi, N. & Yoshimoto, A. eds.), Japan Society of Forest Planning Press, Utsunomiya,
- 吉本敦・柳原宏和・二宮嘉行. 2005. 多変量線形モデルによる林分成長要因探索のための変数選択, *日林誌* 87:504-512

Appendix

$\mathbf{y}_1, \dots, \mathbf{y}_n$ と ζ_1, \dots, ζ_m は独立に多変量正規分布に従う以下のような真のモデルから生成される確率変数とする.

$$[A1] \quad M_* : \mathbf{y}_i \sim i.i.d. N_q(\Xi_*' \mathbf{x}_{d,i}, \Sigma_*) \quad (i = 1, \dots, n)$$

$$M_* : \zeta_j \sim i.i.d. N_q(\Xi_*' \mathbf{w}_{\delta,j}, \Sigma_*) \quad (j = 1, \dots, m)$$

ここで、 $\beta_j = \zeta_j - \Xi_*' \mathbf{w}_{\delta,j}$ と置くと、 $\beta_1, \dots, \beta_m \sim i.i.d. N_q(\mathbf{0}_q, \Sigma_*)$ である. このとき、非線形関数 $f(t | \zeta_j)$ を、テーラー展開を用いた一次近似(参照 Seber and Wild 1989, p.23-25)により近似すると、

$$[A2] \quad f(t | \zeta_j) \approx f(t | \Xi_*' \mathbf{w}_{\delta,j}) + h_j(t | \Xi_*)' \beta_j$$

となる. また、 $krq \times 1$ ベクトル \mathbf{u} を

$$[A3] \quad \mathbf{u} = \sqrt{n} \{ \text{vec}(\hat{\Xi}_Y) - \text{vec}(\Xi_*) \}$$

とすると、 \mathbf{Y} の正規性の仮定より、 \mathbf{u} は以下のような正規分布に従う.

$$[A4] \quad \mathbf{u} \sim N_{krq}(\mathbf{0}_{krq}, \Sigma_* \otimes n(\mathbf{X}_d' \mathbf{X}_d)^{-1})$$

この変数を用いると、テーラー展開により、

$$[A5] \quad f(t | \hat{\Xi}_Y' \mathbf{w}_{\delta,j}) = f(t | \Xi_*' \mathbf{w}_{\delta,j}) + \frac{1}{\sqrt{n}} g_j(t | \Xi_*)' \mathbf{u} + O_p(n^{-1})$$

となる. ここで、 $\sqrt{n/m} = O(1)$ であれば、

$$[A6] \quad \frac{\hat{G}_{cs}(t) - \hat{G}_{cs}(t)}{a_s \sqrt{m}} = \frac{1}{\sqrt{m}} \sum_{j=1}^m h_j(t | \Xi_*)' \beta_j - \sqrt{\frac{m}{n}} \bar{g}(t | \Xi_*)' \mathbf{u} + O_p(n^{-1/2})$$

と展開できる. ただし $a_s = 16.72/61S$ である. ここで β_1, \dots, β_m と \mathbf{u} は互いに独立に正規分布に従う確率変数なので、[A6] 式の右辺は正規分布に収束する. よって、

$$[A7] \quad \frac{\hat{G}_{cs}(t) - \hat{G}_{cs}(t)}{a_s \sqrt{m}} \xrightarrow{D} N(0, \psi^2(t | \Xi_*, \Sigma_*)) \quad (n \rightarrow \infty, \sqrt{n/m} = O(1))$$

を満たす. このとき、 $\hat{\Xi}_Y \rightarrow \Xi_*$ 、 $\hat{\Sigma}_Y \rightarrow \Sigma_*$ ($n \rightarrow \infty$) であるため、 $\psi^2(t | \hat{\Xi}_Y, \hat{\Sigma}_Y) \rightarrow \psi^2(t | \Xi_*, \Sigma_*)$ ($n \rightarrow \infty$) となることがわかる. よって、以下の式が成り立つ.

$$[A8] \quad \frac{\hat{G}_{cs}(t) - \hat{G}_{cs}(t)}{a_s \psi(t | \hat{\Xi}_Y, \hat{\Sigma}_Y) \sqrt{m}} \xrightarrow{D} N(0, 1) \quad (n \rightarrow \infty, \sqrt{n/m} = O(1))$$

上記の式を用いると、

$$[A9] \quad P \left(\frac{|\mathcal{G}_{cs}(t) - \hat{\mathcal{G}}_{cs}(t)|}{a_S \psi(t | \hat{\Xi}_Y, \hat{\Sigma}_Y) \sqrt{m}} \leq z_{\alpha/2} \right) \rightarrow 1 - \alpha \quad (n \rightarrow \infty, \sqrt{n/m} = O(1))$$

となる. [A9] 式を変形することにより [34] 式を得ることができ, 定理を証明することができる.

