

# ニューラルネットワークと回帰分析

021120,030526

## 1 回帰分析と最小2乗解

$$\mathbf{y} = \mathbf{X}\boldsymbol{\beta} + \boldsymbol{\varepsilon} \quad (1)$$

$\mathbf{y} = n \times 1$                     基準変数  
 $\mathbf{X} = n \times (q + 1)$         説明変数 + 1  
 $\boldsymbol{\beta} = (q + 1) \times 1$         回帰係数 (切片を含む)

$$\begin{pmatrix} \hat{y}_1 \\ \hat{y}_{21} \\ \vdots \\ \hat{y}_n \end{pmatrix} = \begin{bmatrix} 1 & x_{11} & x_{12} & \cdots & x_{1q} \\ 1 & x_{21} & x_{22} & \cdots & x_{2q} \\ \vdots & \vdots & \vdots & \cdots & \vdots \\ 1 & x_{n1} & x_{n2} & \cdots & x_{nq} \end{bmatrix} \times \begin{pmatrix} \mu \\ \beta_1 \\ \beta_2 \\ \vdots \\ \beta_q \end{pmatrix} \quad (2)$$

$$\text{RSS} = \sum_{i=1}^n (y_i - \hat{y}_i)^2 \quad (3)$$

$$\hat{y}_i = \sum_{j=0}^q x_{ij}\beta_j \quad (4)$$

$$\frac{\partial \text{RSS}}{\partial \beta_k} = 2 \sum_{i=1}^n (y_i - \hat{y}_i) \times \frac{\partial (y_i - \hat{y}_i)}{\partial \beta_k} \quad (5)$$

$$= 2 \sum_{i=1}^n (y_i - \hat{y}_i) \times \frac{\partial \sum_{j=0}^q x_{ij}\beta_j}{\partial \beta_k} \quad (6)$$

$$= 2 \sum_{i=1}^n (y_i - \hat{y}_i) \times x_{ik} \quad (7)$$

$$= 2 \sum_{i=1}^n (y_i - \sum_{j=0}^q x_{ij}\beta_j) \times x_{ik} \quad (8)$$

$$= 2 \sum_{i=1}^n y_i x_{ik} - 2 \sum_{i=1}^n \left( \sum_{j=0}^q x_{ij}\beta_j \right) x_{ik} \quad (9)$$

$$= 2 \sum_{i=1}^n y_i x_{ik} - 2 \sum_{i=1}^n \left( \sum_{j=0}^q x_{ij} x_{ik} \beta_j \right) \quad (10)$$

$$= 2 \sum_{i=1}^n y_i x_{ik} - 2 \sum_{j=0}^q \left( \sum_{i=1}^n x_{ij} x_{ik} \right) \beta_j \quad (11)$$

$\sum_{i=1}^n y_i x_{ik}$  は  $\mathbf{X}'\mathbf{y}$  の第  $k$  要素、 $\sum_{i=1}^n x_{ij} x_{ik}$  は  $\mathbf{X}'\mathbf{X}$  の第  $jk$  要素  
 したがって、方程式  $\frac{\partial \text{RSS}}{\partial \beta_k} = 0, k = 0, 1, 2, \dots, q$  は、まとめて

$$\mathbf{X}'\mathbf{y} - \mathbf{X}'\mathbf{X}\boldsymbol{\beta} = \mathbf{0} \quad (12)$$

と書ける。従って、解は

$$\hat{\beta} = (\mathbf{X}'\mathbf{X})^{-1}\mathbf{X}'\mathbf{y} \quad (13)$$

となる。

## 2 バリエーション

### 1. 単回帰分析

基準変数 1 個、説明変数 1 個

$$\mathbf{y} = \beta_0\mathbf{1} + \beta\mathbf{x} + \varepsilon \quad (14)$$

### 2. 重回帰分析

基準変数 1 個、説明変数  $q$  個

$$\mathbf{y} = \mathbf{X}\beta + \varepsilon \quad (15)$$

### 3. 多変量回帰分析

基準変数  $p$  個、説明変数  $q$  個

$$\mathbf{Y} = \mathbf{X}\beta' + \varepsilon, \quad \beta = (p \times q) \quad (16)$$

$$\hat{\beta}' = (\mathbf{X}'\mathbf{X})^{-1}\mathbf{X}'\mathbf{Y} \quad (17)$$

### 4. 冗長性分析

$$\mathbf{Y} = \mathbf{1}\mu' + \mathbf{X}\mathbf{A}\mathbf{B}' + \varepsilon, \quad \mathbf{A} = (q \times r), \quad \mathbf{B} = (p \times r) \quad r < p, q \quad (18)$$

### 5. 多項式回帰分析

$$\hat{y}_i = \sum_{m=0}^M \beta_m x_i^m \quad (19)$$

$$\mathbf{Y} = \mathbf{G}^{(M)}(\mathbf{x})\beta + \varepsilon, \quad \mathbf{G}^{(M)}(\mathbf{x}) = (n \times M + 1) \quad (20)$$

$$\mathbf{G}(\mathbf{x}) = [\mathbf{g}_m(\mathbf{x})], \quad m = 0, 1, \dots, M \quad (21)$$

$$\mathbf{g}_m(\mathbf{x}) = [x_i^m] \quad (22)$$

### 6. スプライン回帰分析 (多項式 + 切断冪関数)

$$\hat{y}_i = \sum_{m=0}^M \beta_m x_i^m + \sum_{\ell=1}^K \beta_\ell (x_i - c_\ell)_+^M \quad (23)$$

$$(x_i - c_\ell)_+^M = \begin{cases} 0 & (x_i < c_\ell) \\ (x_i - c_\ell)^M & (x_i \geq c_\ell) \end{cases} \quad (24)$$

$$\mathbf{Y} = \mathbf{G}^{(K,M)}(\mathbf{x})\beta + \varepsilon \quad \mathbf{G}^{(K,M)}(\mathbf{x}) = (n \times M + K + 1) \quad (25)$$

ただし、 $K$  は節点の数、 $c_\ell$  は節点

## 7. ニューラルネットワーク

$$y_i = \beta_0 + \sum_{j=1}^Q \psi(a_j(x_i - b_j))\beta_j + \varepsilon_i, \quad i = 1, 2, \dots, n \quad (26)$$

ただし、 $a_j, b_j$  は定数であり、

$$\psi(x) = \frac{1}{(1 + \exp(-x))} \quad (27)$$

はロジスティック関数である。従って

$$\mathbf{y} = \mathbf{G}(\mathbf{x})\boldsymbol{\beta} + \boldsymbol{\varepsilon} \quad (28)$$

の形で書くと

$$\mathbf{G}(\mathbf{x}) = [\mathbf{1}, \psi(a_1(\mathbf{x} - b_1)), \psi(a_2(\mathbf{x} - b_2)), \dots, \psi(a_Q(\mathbf{x} - b_Q))] \quad (29)$$

となる。

なお、ニューラルネットワークモデルでは、 $w_{ij}^{(k)}$  を第  $k-1$  層の第  $i$  番目の変数から第  $k$  層の第  $j$  番目の変数を作るための重み、また、その際用いられる定数項を  $\theta_j^{(k)}$  として、以下のように記すことが多い。

$$y_i = \theta^{(3)} + \sum_{j=1}^Q \psi(w_{j1}^{(2)} x_i + \theta_j^{(2)}) w_{j1}^{(3)} + \varepsilon_i, \quad i = 1, 2, \dots, n \quad (30)$$

上記のモデルは、 $x$  を第 1 層の一つしかないニューロンの出力とし、第 2 層の第  $j$  番目のニューロンの入力が  $w_{j1}^{(2)} x_i + \theta_j^{(2)}$ 、その出力が  $\psi(w_{j1}^{(2)} x_i + \theta_j^{(2)})$  となっていると考えられる。したがって、これまでの係数が既知の場合は、ニューラルネットワークモデルは、 $n_2$  個の第 2 層の出力を説明変数、 $y$  を基準変数とした重回帰分析モデルと考えることが出来る。

スプライン回帰や多項式と似ている点は、両者とも、本来の説明変数を  $G$  という形に展開し、それらの線形結合で  $\mathbf{y}$  を予測しようと試みることである。

両者の違いは、スプラインの場合は、展開に際して用いられる接点  $c_\ell, \ell = 1, 2, \dots, K$  は定数であるのに対し、ニューラルネットワークにおいては、その係数  $(a_j, b_j), j = 1, 2, \dots, Q$  はパラメタとして推定されるものである点である。

また、多項式回帰の場合、中間層を増やすことは、結局、最初に展開された説明変数の線形回帰を作ることになり、意味がないが、スプラインや、特にニューラルネットワークの場合は、意味がある。

## 3 ニューラルネットワークについて

$Q$  個の説明変数  $\mathbf{x}$  を第 1 層、 $Q$  個の基準変数  $\mathbf{y}$  を第  $n$  層とするネットワークを考える。各層のニューロンの数は  $n_k, k = 1, 2, \dots, n$  とする。従って、 $n_1 = Q, n_n = Q$  である。また、第  $k$  層の第  $j$  番目のニューロンへの入力を  $i_j^k$ 、同じく出力を  $o_j^k$  と書く。各個体を表す添え字は  $a, a = 1, 2, \dots, N$  である。各ニューロンはそれぞれ上位のニューロンに向かって同一の出力を発する。また、各ニューロンは下位のニューロンの発する出力を重み付けした形でそのニューロンへの入力形成する。すなわち、第  $k-1$  層の第  $i$  番目のニューロンの出力  $o_i^{k-1}$  から第  $k$  層の第  $j$  番目のニューロンの入力  $i_j^k$  を作る際の重みを  $w_{ij}^k$  とすれば、

$$i_{aj}^k = \sum_{i=1}^{n_{k-1}} o_{ai}^{k-1} w_{ij}^k + \theta_j^k \quad (31)$$

である。また、第  $k$  層の入力はシグモイド関数  $\psi$  により、

$$o_{aj}^k = \psi(i_{aj}^k) = \frac{1}{1 + \exp(-\lambda i_{aj}^k)}, \quad (\lambda = 1.7 \text{ or } 1) \quad (32)$$

という変換を受けた形で出力される。

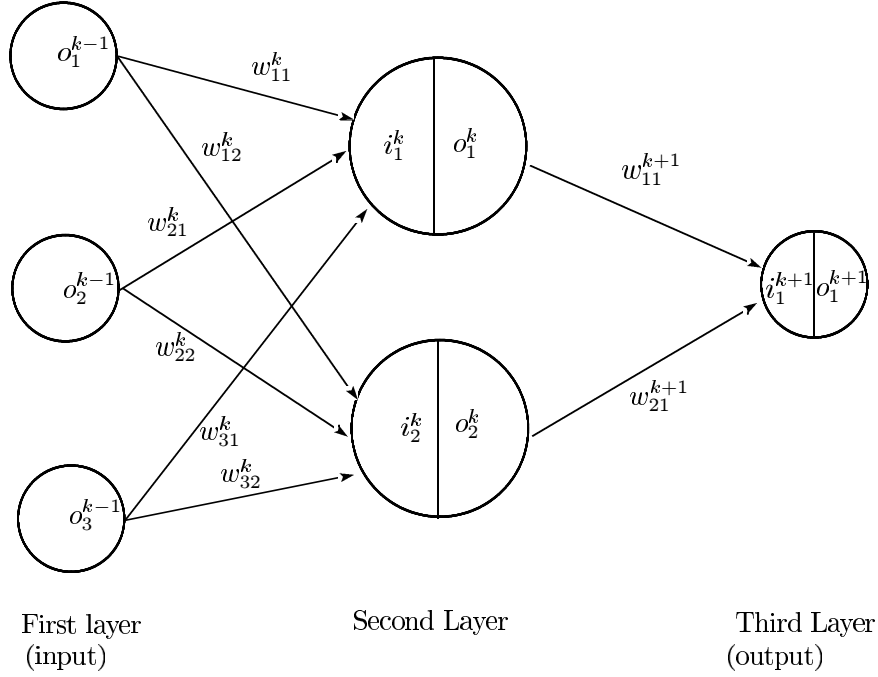


Figure 1: ニューラルネットワークモデル

これを全ての  $i_{aj}^k, o_{aj}^k, a = 1, 2, \dots, N$  に関してまとめて書くと、

$$\mathbf{O}^k = \Psi(\mathbf{I}^k) \quad (33)$$

$$= \Psi(\mathbf{O}^{k-1} \mathbf{W}^k + \mathbf{1}\theta^k) \quad (34)$$

$$= \Psi(\Psi(\mathbf{I}^{k-1}) \mathbf{W}^k + \mathbf{1}\theta^k) \quad (35)$$

$$= \Psi(\Psi(\mathbf{O}^{k-2} \mathbf{W}^{k-1} + \mathbf{1}\theta^{k-1}) \mathbf{W}^k + \mathbf{1}\theta^k) \quad (36)$$

$$\vdots \quad (37)$$

となるが、 $(\mathbf{O}^k : N \times n_k, \mathbf{W}^k : n_{k-1} \times n_k, \theta_j : 1 \times n_k)$  これは

$$\mathbf{O}^k = \Psi(\mathbf{O}^{k-1} \mathbf{W}^k + \mathbf{1}\theta^k) = \Gamma_k(\mathbf{O}^{k-1}) \quad (38)$$

と定義をすれば、最終層の出力が

$$\mathbf{O}^n = \Gamma_n(\mathbf{O}^{n-1}) \quad (39)$$

$$= \Gamma_n(\Gamma_{n-1}(\mathbf{O}^{n-2})) \quad (40)$$

$$\vdots \quad (41)$$

$$= \Gamma_n \circ \Gamma_{n-1} \circ \dots \circ \Gamma_2(\mathbf{O}^1) \quad (42)$$

$$= \Gamma(\mathbf{O}^1) \quad (43)$$

と書けることを示している。ただし、

$$\Gamma(\cdot) = \Gamma_n \circ \Gamma_{n-1} \circ \dots \circ \Gamma_2(\cdot) \quad (44)$$

である。従って、ニューラルネットワークのモデルは、 $\mathbf{Y} = \mathbf{O}$  とおくことにより、

$$\mathbf{Y} = \Gamma(\mathbf{X}) + \boldsymbol{\varepsilon} \quad \text{or} \quad \hat{\mathbf{Y}} = \Gamma(\mathbf{X}) \quad (45)$$

と書くことができる。

## 4 最小2乗法によるパラメタの推定

データ  $\mathbf{Y}$  とモデルの値  $\hat{\mathbf{Y}}$  の差の2乗和

$$\text{RSS} = \sum_{a=1}^N \sum_{j=1}^Q (y_{aj} - \hat{y}_{aj})^2 \quad (46)$$

を最小にするようにモデルのパラメタ

$$w_{ij}^k, \theta_j^k, k = 2, 3, \dots, n, i = 1, 2, \dots, n_{k-1}, j = 1, 2, \dots, n_k$$

を求めるためには以下に示す  $\hat{y}_{aj}$  を各パラメタで偏微分したものが必要となる。ただし、

$$\hat{y}_{aj} = o_{as}^n \quad (47)$$

であるが、個体を表す添字  $a$  は落としてある。また、 $\theta$  に関する微分も同時に示してある。

$$\frac{\partial o_s^n}{\partial w_{ij}^k} = \frac{\partial o_s^n}{\partial i_j^k} \times \frac{\partial i_j^k}{\partial w_{ij}^k} \quad (48)$$

$$\frac{\partial o_s^n}{\partial i_j^k} = \frac{\partial o_s^n}{\partial o_j^k} \times \frac{\partial o_j^k}{\partial i_j^k} \quad \text{and} \quad \frac{\partial o_s^n}{\partial o_j^k} = \sum_{q=1}^{n_{k+1}} \frac{\partial o_s^n}{\partial i_q^{k+1}} \times \frac{\partial i_q^{k+1}}{\partial o_j^k} \quad (49)$$

$$\frac{\partial i_q^{k+1}}{\partial o_j^k} = w_{jq}^{k+1} \quad \text{and} \quad \frac{\partial o_j^k}{\partial i_j^k} = o_j^k(1 - o_j^k) \quad (50)$$

従って、 $\delta_j^{k,s} = \frac{\partial o_s^n}{\partial i_j^k}$  と書くと、

$$\delta_j^{k,s} = \frac{\partial o_s^n}{\partial i_j^k} \quad (51)$$

$$= \sum_{q=1}^{n_{k+1}} \frac{\partial o_s^n}{\partial i_q^{k+1}} w_{jq}^{k+1} \times \frac{\partial o_j^k}{\partial i_j^k} = \sum_{q=1}^{n_{k+1}} \delta_q^{k+1,s} w_{jq}^{k+1} \times \frac{\partial o_j^k}{\partial i_j^k} \quad (52)$$

と再帰的にかける。(これが誤差逆伝搬の由来である。) ただし

$$\frac{\partial i_j^k}{\partial w_{ij}^k} = \begin{cases} o_i^{k-1} & \text{if not } \theta \\ 1 & \text{if } \theta \end{cases} \quad (53)$$

である。また、 $k = n$  の場合には

$$\delta_j^{n,s} = \frac{\partial o_s^n}{\partial i_j^n} = \begin{cases} o_s^n(1 - o_s^n) & \text{if } j = s \\ 0 & \text{otherwise} \end{cases} \quad (54)$$

となる。

最適化の具体的方法に関しては次節以降に述べるが、ニューラルネットワークの重みの推定には、全ての  $\mathbf{W}^K$  もしくはその部分を  $R \times 1$  のベクトルの形にしたものを  $\boldsymbol{\xi}$  とし、ヤコビアン行列  $\mathbf{F}_{(j)}(\boldsymbol{\xi})$  の要素としてこの  $\frac{\partial o_j^n}{\partial w_{IJ}^k}$  を用いることになる。また、個体を表す添字  $a$  は、次節では  $i$  となっている。

## 5 非線形モデルの最小2乗解： ガウス・ニュートン法

モデル  $\hat{y}(\boldsymbol{\xi})$  がパラメタ  $\boldsymbol{\xi}$  ( $R \times 1$ ) に関して非線形の場合、最小2乗基準

$$\text{RSS}(\boldsymbol{\xi}) = \|\mathbf{y} - \hat{\mathbf{y}}(\boldsymbol{\xi})\|^2 = \sum_{i=1}^N (y_i - \hat{y}_i(\boldsymbol{\xi}))^2 \quad (55)$$

を最小にするようなパラメタの推定値を解析的に求めることは困難である場合が多い。しかし、そのような場合にも、何らかの数値的最適化の方法を用いれば、最小2乗解を求めることが可能である。

一般に、ある関数  $r(\boldsymbol{\xi})$  の極値（極小値または極大値のことであるが、本節では極小値とする）を与えるパラメタ  $\boldsymbol{\xi}$  の値を数値的に求めることを関数の最適化と呼ぶが、その方法は、何らかの方法で求めた  $\boldsymbol{\xi}$  の初期値から、 $r(\boldsymbol{\xi})$  が減少する方向  $\mathbf{d}$  へある分量  $\alpha$  だけ修正する、

$$\boldsymbol{\xi}_{(\ell+1)} = \boldsymbol{\xi}_{(\ell)} + \alpha \mathbf{d}_{(\ell)} \quad (56)$$

という形の繰り返し計算でパラメタの値を更新していくという形をとる。数値的最適化の方法には様々なものがあるが、それらの主な違いは、更新の方向  $\mathbf{d}_{(\ell)}$  をどの様にして定めるかによる。

最も単純なものは最急降下法と呼ばれる方法であり、修正の方向を  $r(\boldsymbol{\xi})$  のパラメタの各要素に関する一次偏微分係数行列の逆方向

$$\mathbf{d}_{(\ell)} = -\mathbf{g}(\boldsymbol{\xi}_{(\ell)}) = - \left. \frac{\partial r(\boldsymbol{\xi})}{\partial \boldsymbol{\xi}} \right|_{\boldsymbol{\xi}=\boldsymbol{\xi}_{(\ell)}} \quad (57)$$

に取るものである。この、 $R \times 1$  のベクトル  $\mathbf{g}$  はグレイディエントと呼ばれることがあるが、これは  $\boldsymbol{\xi}_{(\ell)}$  における  $r(\boldsymbol{\xi})$  のパラメタの各要素に関する傾きの逆方向を示すもので、Eq.(56) を用いることにより、

$$r(\boldsymbol{\xi}_{(\ell+1)}) \leq r(\boldsymbol{\xi}_{(\ell)}) \quad (58)$$

を満たす  $\alpha$  が存在することが知られている。

$r(\boldsymbol{\xi})$  を  $\boldsymbol{\xi}_{(\ell)}$  のまわりで2次の項まで Taylor 展開すると

$$r(\boldsymbol{\xi}) \approx r(\boldsymbol{\xi}_{(\ell)}) + \mathbf{g}(\boldsymbol{\xi}_{(\ell)})'(\boldsymbol{\xi} - \boldsymbol{\xi}_{(\ell)}) + \frac{1}{2}(\boldsymbol{\xi} - \boldsymbol{\xi}_{(\ell)})' \mathbf{H}(\boldsymbol{\xi}_{(\ell)})(\boldsymbol{\xi} - \boldsymbol{\xi}_{(\ell)}) \quad (59)$$

となるが、この  $R \times R$  の2次微分行列（ $r$  をパラメタの各要素で2回偏微分したものを要素として持つ行列）

$$\mathbf{H}(\boldsymbol{\xi}_{(\ell)}) = \left[ \left. \frac{\partial^2 r(\boldsymbol{\xi})}{\partial \boldsymbol{\xi} \partial \boldsymbol{\xi}'} \right|_{\boldsymbol{\xi}=\boldsymbol{\xi}_{(\ell)}} \right] \quad (60)$$

はヘシアン行列と呼ばれるものである。これを考慮すると、Eq.(56) において

$$\mathbf{d}_{(\ell)} = -\mathbf{H}(\boldsymbol{\xi}_{(\ell)})^{-1} \mathbf{g}(\boldsymbol{\xi}_{(\ell)}) \quad (61)$$

とおいたものが Newton-Raphson 法の更新式として得られる。

最適化されるべき関数が最小2乗基準の場合（ $r(\boldsymbol{\xi}) = \text{RSS}(\boldsymbol{\xi})$ ）求めるべきパラメタの仮の値を  $\boldsymbol{\xi}_{(\ell)}$  とする。モデルの各要素を  $\boldsymbol{\xi}_{(\ell)}$  のまわりで1次の項まで Taylor 展開したものをまとめれば、

$$\hat{\mathbf{y}}(\boldsymbol{\xi}) \approx \hat{\mathbf{y}}(\boldsymbol{\xi}_{(\ell)}) + \mathbf{F}(\boldsymbol{\xi}_{(\ell)})(\boldsymbol{\xi} - \boldsymbol{\xi}_{(\ell)}) \quad (62)$$

となるが、ここに、

$$\mathbf{F}(\boldsymbol{\xi}) = \frac{\partial \hat{\mathbf{y}}(\boldsymbol{\xi})}{\partial \boldsymbol{\xi}} = \left[ \frac{\partial \hat{y}_i(\boldsymbol{\xi})}{\partial \xi_k} \right], \quad i = 1, 2, \dots, N, \quad k = 1, 2, \dots, R \quad (63)$$

は、ヤコビアン行列と呼ばれる  $N \times R$  の 1 次微分行列、また、 $\mathbf{F}(\boldsymbol{\xi}_{(\ell)})$  はそれを  $\boldsymbol{\xi}_{(\ell)}$  で評価した行列

$$\mathbf{F}(\boldsymbol{\xi}_{(\ell)}) = \left. \frac{\partial \hat{\mathbf{y}}(\boldsymbol{\xi})}{\partial \boldsymbol{\xi}} \right|_{\boldsymbol{\xi}=\boldsymbol{\xi}_{(\ell)}} \quad (64)$$

である。このことを考慮して、Eq.(56) において

$$\mathbf{d}_{(\ell)} = -\frac{1}{2}(\mathbf{F}(\boldsymbol{\xi}_{(\ell)})' \mathbf{F}(\boldsymbol{\xi}_{(\ell)}))^{-1} \mathbf{g}(\boldsymbol{\xi}_{(\ell)}) \quad (65)$$

と置いたものを用いる方法が Gauss-Newton 法と呼ばれるものである。なお、このヤコビアン行列を用いると

$$\mathbf{g} = \frac{\partial \text{RSS}(\boldsymbol{\xi})}{\partial \boldsymbol{\xi}} = -2\mathbf{F}(\boldsymbol{\xi})'(\mathbf{y} - \hat{\mathbf{y}}(\boldsymbol{\xi})) \quad (66)$$

となる。Gauss-Newton 法は 1 次の情報しか利用しないが効率はよく、Eq.(58) を満たす  $\alpha$  も必ず存在することが知られている。

また、基準変数が多変量の場合は、

$$\text{RSS} = \|\mathbf{y} - \hat{\mathbf{y}}(\boldsymbol{\xi})\|^2 = \sum_{j=1}^p \|\mathbf{y}_{(j)} - \hat{\mathbf{y}}_{(j)}(\boldsymbol{\xi})\|^2 = \sum_{j=1}^p \sum_{i=1}^N (y_{ij} - \hat{y}_{ij}(\boldsymbol{\xi}))^2 \quad (67)$$

と書けるため、変数ごとのヤコビアン行列を

$$\mathbf{F}_{(j)}(\boldsymbol{\xi}) = \frac{\partial \hat{\mathbf{y}}_{(j)}(\boldsymbol{\xi})}{\partial \boldsymbol{\xi}} = \left[ \frac{\partial \hat{y}_{ij}(\boldsymbol{\xi})}{\partial \boldsymbol{\xi}_k} \right], \quad i = 1, 2, \dots, N, \quad k = 1, 2, \dots, R \quad (68)$$

と定義すれば、グレイディエント並びに更新の方向は

$$\mathbf{g} = \frac{\partial \text{RSS}(\boldsymbol{\xi})}{\partial \boldsymbol{\xi}} = -2 \sum_{j=1}^p \mathbf{F}_{(j)}(\boldsymbol{\xi})'(\mathbf{y}_{(j)} - \hat{\mathbf{y}}_{(j)}(\boldsymbol{\xi})) \quad (69)$$

$$\mathbf{d}_{(\ell)} = -\frac{1}{2} \left( \sum_{j=1}^p \mathbf{F}_{(j)}(\boldsymbol{\xi}_{(\ell)})' \mathbf{F}_{(j)}(\boldsymbol{\xi}_{(\ell)}) \right)^{-1} \mathbf{g}(\boldsymbol{\xi}_{(\ell)}) \quad (70)$$

で与えられる。

パラメタを更新すべき方向が定めれば、ステップサイズ  $\alpha$  の値を適当に決めることにより、Eq.(58) を満たすような  $\boldsymbol{\xi}_{(\ell)}$  を求めることが出来る。ステップサイズを決める最も簡単な方法は、まず  $\alpha = 1$  から始めて Eq.(58) が満たされない場合には、それが満たされるまで  $\alpha$  の値を半減していくという方法である。より高度な方法としては、RSS の変化を参考にしてステップサイズの減少幅を決める方法や、 $\mathbf{d}_{(\ell)}$  を既知として RSS を  $\alpha$  の関数とみなし、それを 3 次関数で近似して極値を求める方法がある。

Gauss-Newton 法の一般的な形は以下のようなになる。

### Gauss-Newton 法

$\boldsymbol{\xi}_1$  に初期値を設定する。

for  $\ell = 1$  to maxiter until(  $\text{RSS}(\boldsymbol{\xi}_{(\ell)}) - \text{RSS}(\boldsymbol{\xi}_{(\ell+1)})$  が小さい );

$\mathbf{F}_{(j)}(\boldsymbol{\xi}_{(\ell)})$ ,  $\mathbf{g}(\boldsymbol{\xi}_{(\ell)})$ ,  $\mathbf{d}_{(\ell)}$  を計算する。

$\alpha=1$ ;

for iters = 1 to maxhalf;

$\boldsymbol{\xi}_{(\ell+1)} = \boldsymbol{\xi}_{(\ell)} + \alpha \mathbf{d}_{(\ell)}$ ;

if  $\text{RSS}(\boldsymbol{\xi}_{(\ell+1)}) \leq \text{RSS}(\boldsymbol{\xi}_{(\ell)})$  then break;

$\alpha=\alpha/2$ ;

end;

end;