

# 統計モデルのパラメタの推定

前川眞一

東京工業大学 大学院社会理工学研究科

## 1 統計モデル

$N$  個の個体に関するデータ  $(\mathbf{y}_i, \mathbf{x}_i), i = 1, 2, \dots, N$  が与えられた場合、統計学では以下のようなモデルを用いてこのデータの分析を行うことが多い。

$$\mathbf{y}_i = f(\mathbf{x}_i|\boldsymbol{\beta}) + \varepsilon_i, i = 1, 2, \dots, N \quad (1)$$

ただし、 $N$  は観測値（個体）の数、 $\mathbf{x}_i$  は  $q \times 1$  の独立変数（説明変数）、 $\mathbf{y}_i$  は  $p \times 1$  の従属変数（基準変数）、 $\varepsilon_i$  は  $p \times 1$  の誤差項（残差）、また、 $\boldsymbol{\beta}$  はパラメタベクトルである。この形は、観測されるデータ  $\mathbf{y}$  から誤差を取り除いた部分は、 $\mathbf{x}$  の関数  $f$  として記述できることを示しているが、これは、 $\mathbf{x}$  の値を変化させれば  $\mathbf{y}$  の値も変化し、その両者の関係が  $f$  という関数で記述出来ることを意味している。また、通常は、 $f$  には自由に決めることの出来るパラメタ  $\boldsymbol{\beta}$  が含まれている場合が多い。このモデルは  $\mathbf{y}$  から誤差を取り除いた部分、すなわち、システムティックに変化する部分を  $\hat{\mathbf{y}}$  とし、それを  $\mathbf{y}$  の  $\mathbf{x}$  を与えたときの条件付き期待値と考えれば、

$$\mathbf{y}_i = \hat{\mathbf{y}}_i + \varepsilon_i \quad \text{and} \quad \mathcal{E}(\mathbf{y}_i | \mathbf{x}_i) = \hat{\mathbf{y}}_i \quad \text{or} \quad \mathcal{E}(\varepsilon_i) = 0 \quad (2)$$

$$\hat{\mathbf{y}}_i = f(\mathbf{x}_i|\boldsymbol{\beta}) \quad (3)$$

と書くこともできる。また、全データに関してこの関係が成り立つことを、

$$\mathbf{Y} = F(\mathbf{X}|\boldsymbol{\beta}) + \boldsymbol{\varepsilon} \quad \text{or} \quad \mathbf{Y} = \hat{\mathbf{Y}} + \boldsymbol{\varepsilon} \quad (4)$$

とも書くことにする。ただし、 $\mathbf{Y}$  と  $\hat{\mathbf{Y}}$  は  $N \times p$ 、 $\mathbf{X}$  は  $N \times q$ 、 $\boldsymbol{\varepsilon}$  は  $N \times p$ 、 $F$  は  $N \times p$  の行列である。なお、従属変数の数が  $p = 1$  の場合を1変数（1変量）モデル、それ以外を多変量モデルと呼ぶことがある。

例えば、この形のモデルの最も単純なものは、分散分析や線形回帰分析のモデルであり、非線形の回帰分析（例えばニューラルネットワークモデル）や因子分析、主成分分析、クラスター分析、多次元尺度法等のモデルもこの形で説明することが出来る。

本稿においては、データが与えられた場合に、上記のモデルのパラメタを推定する方法を説明する。なお、パラメタの推定方法は様々な方式があるが、ここでは一般的に利用が可能な最尤法と最小2乗法に付いての説明を行う。

なお、データ  $\mathbf{y}_i$  がカテゴリカルな場合には、その期待値  $\hat{\mathbf{y}}_i$  は生起確率となり、通常は2項分布や多項分布を仮定した最尤解を用いることになる。

## 2 データの種類と適合度の基準

### 2.1 計量的データ

基準変数  $\mathbf{Y}$  が計量データの場合、誤差項の分布として正規分布を仮定する場合が多い。

$$\varepsilon_i \sim N(0, \sigma^2) \quad \text{or} \quad \varepsilon_i \sim N(\mathbf{0}, \boldsymbol{\Sigma}) \quad (5)$$

この場合は  $Y$  の  $X$  を与えたときの条件付き密度関数  $g$  が

$$g(y_i | \mathbf{X}, \boldsymbol{\beta}) = \frac{1}{\sqrt{2\pi\sigma^2}} \exp\left(-\frac{1}{2}(y_i - f(\mathbf{x}_i | \boldsymbol{\beta}))^2 / \sigma^2\right) \quad (6)$$

$$g(\mathbf{y}_i | \mathbf{X}, \boldsymbol{\beta}) = \frac{1}{\sqrt{2\pi|\boldsymbol{\Sigma}|}} \exp\left(-\frac{1}{2}(\mathbf{y}_i - f(\mathbf{x}_i | \boldsymbol{\beta}))' \boldsymbol{\Sigma}^{-1}(\mathbf{y}_i - f(\mathbf{x}_i | \boldsymbol{\beta}))\right) \quad (7)$$

と書けるため、 $N$  個の独立な観測値を得た場合、全てのデータの同時密度は

$$g(\mathbf{y} | \mathbf{X}, \boldsymbol{\beta}) = \prod_{i=1}^N \frac{1}{\sqrt{2\pi\sigma^2}} \exp\left(-\frac{1}{2}(y_i - f(\mathbf{x}_i | \boldsymbol{\beta}))^2 / \sigma^2\right) \quad (8)$$

$$g(\mathbf{y} | \mathbf{X}, \boldsymbol{\beta}) = \prod_{i=1}^N \frac{1}{\sqrt{2\pi|\boldsymbol{\Sigma}|}} \exp\left(-\frac{1}{2}(\mathbf{y}_i - f(\mathbf{x}_i | \boldsymbol{\beta}))' \boldsymbol{\Sigma}^{-1}(\mathbf{y}_i - f(\mathbf{x}_i | \boldsymbol{\beta}))\right) \quad (9)$$

となる。このデータの同時密度関数をパラメタ  $\boldsymbol{\beta}$  の関数と見たものを尤度関数  $L$ 、また、その対数を対数尤度関数  $\ln L$  と呼ぶ。

$$L(\boldsymbol{\beta}) = \prod_{i=1}^N \frac{1}{\sqrt{2\pi\sigma^2}} \exp\left(-\frac{1}{2}(y_i - f(\mathbf{x}_i | \boldsymbol{\beta}))^2 / \sigma^2\right) \quad (10)$$

$$\ln L(\boldsymbol{\beta}) = -\frac{1}{2} \left( \ln 2\pi\sigma^2 + \sum_{i=1}^N (y_i - f(\mathbf{x}_i | \boldsymbol{\beta}))^2 / \sigma^2 \right) \quad (11)$$

この尤度を最大とするようなパラメタの値をパラメタの最尤解と呼ぶ。また、対数尤度の形から、誤差分散が全て等しい場合には、この解は

$$\text{RSS} = \sum_{i=1}^N (y_i - f(\mathbf{x}_i | \boldsymbol{\beta}))^2 \quad (12)$$

$$\text{RSS} = \sum_{i=1}^N (\mathbf{y}_i - f(\mathbf{x}_i | \boldsymbol{\beta}))' \boldsymbol{\Sigma}^{-1} (\mathbf{y}_i - f(\mathbf{x}_i | \boldsymbol{\beta})) \quad (13)$$

を最少とする最小2乗解と同じであると言える。(ただし、 $\boldsymbol{\Sigma}^{-1}$  行列は既知の重み行列として取り扱う。) なお、多変量の場合で重み行列がない場合は、

$$\text{RSS} = \sum_{j=1}^p \sum_{i=1}^N (y_{ij} - f_j(\mathbf{x}_i | \boldsymbol{\beta}))^2 \quad (14)$$

$$(15)$$

と書くことができる。

### 2.1.1 例

計量データ：正規分布

- 平均値の推定：  $\mathbf{X}$  は全て 1
- 多群の平均値の推定：  $\mathbf{X}$  は群を示すダミー変数（デザイン行列）
- 回帰分析：  $\hat{\mathbf{Y}} = \mathbf{X}\boldsymbol{\beta}'$
- ニューラルネット：  $\hat{\mathbf{Y}} = \boldsymbol{\beta}$  の要素の非線形関数
- 因子分析：  $\mathbf{X}$  は計量的潜在変数（これも推定する）
- クラスター分析：  $\mathbf{X}$  は計数的潜在変数（これも推定する）
- 多次元尺度法：  $\mathbf{X}$  は刺激対を示すダミー変数（デザイン行列）

## 2.2 計数的データ（度数データ、カテゴリカルデータ）

データが、ある個体が  $p$  個のカテゴリの中でどのカテゴリに反応したか、というような形で得られる場合、それを計数的データ、度数データ、もしくはカテゴリカルデータと呼ぶ。この場合、 $\mathbf{y}_i$  はその個体がどのカテゴリに反応したかを示すインディケータベクトルとなる。また、 $\mathbf{y}_i$  の期待値  $\hat{\mathbf{y}}_i$  はその生起確立となる。以下に、そのことを明記するため、生起確立を

$$\hat{\mathbf{y}}_i = f(\mathbf{x}_i|\boldsymbol{\beta}) = \pi(\mathbf{x}_i|\boldsymbol{\beta}) \quad (16)$$

と記す。

基準変数  $\mathbf{Y}$  が計数データの場合、その分布としてベルヌーイ分布を仮定する場合が多い。この場合は  $\mathbf{Y}$  の  $\mathbf{X}$  を与えたときの条件付き確率  $g$  が

$$g(\mathbf{y}_i|\mathbf{X}, \boldsymbol{\beta}) = \prod_{j=1}^p \pi_j(\mathbf{x}_i|\boldsymbol{\beta})^{y_{ij}}, \quad \sum_{j=1}^p \pi_j(\mathbf{x}_i|\boldsymbol{\beta}) = 1 \quad (17)$$

と書けるため、 $N$  個の独立な観測値を得た場合、全てのデータの同時密度や対数尤度関数は

$$g(\mathbf{y}|\mathbf{X}, \boldsymbol{\beta}) = \prod_{i=1}^N \prod_{j=1}^p \pi_j(\mathbf{x}_i|\boldsymbol{\beta})^{y_{ij}} \quad (18)$$

$$\ln L(\boldsymbol{\beta}) = \sum_{i=1}^N \sum_{j=1}^p y_{ij} \ln \pi_j(\mathbf{x}_i|\boldsymbol{\beta}) + \boldsymbol{\beta} \text{ を含まない部分} \quad (19)$$

となる。

なお、この形とは少し異なるが、データ  $\mathbf{y}$  が、 $p$  個の2値変数からなる場合、条件付きの独立性を仮定して、

$$g(\mathbf{y}_i|\mathbf{X}, \boldsymbol{\beta}) = \prod_{j=1}^p \pi_j(\mathbf{x}_i|\boldsymbol{\beta})^{y_{ij}} (1 - \pi_j(\mathbf{x}_i|\boldsymbol{\beta}))^{1-y_{ij}} \quad (20)$$

また、 $(n_{ij}, f_{ij})$ ,  $i = 1, 2, \dots, j = 1, 2, \dots, p$  を各個体の変数ごとの試行数と成功の度数として、度数の条件付き確率が

$$g(\mathbf{f}_i|\mathbf{X}, \boldsymbol{\beta}) \propto \prod_{j=1}^p \pi_j(\mathbf{x}_i|\boldsymbol{\beta})^{f_{ij}} (1 - \pi_j(\mathbf{x}_i|\boldsymbol{\beta}))^{n_{ij}-f_{ij}} \quad (21)$$

と書き表されるモデルを考えることもある。この場合、全てのデータの同時密度は

$$g(\mathbf{y}|\mathbf{X}, \boldsymbol{\beta}) = \prod_{i=1}^N \prod_{j=1}^p \pi_j(\mathbf{x}_i|\boldsymbol{\beta})^{y_{ij}} (1 - \pi_j(\mathbf{x}_i|\boldsymbol{\beta}))^{1-y_{ij}} \quad (22)$$

$$g(\mathbf{f}|\mathbf{X}, \boldsymbol{\beta}) = \prod_{i=1}^N \prod_{j=1}^p \pi_j(\mathbf{x}_i|\boldsymbol{\beta})^{f_{ij}} (1 - \pi_j(\mathbf{x}_i|\boldsymbol{\beta}))^{n_{ij}-f_{ij}} \quad (23)$$

となり、パラメタの対数尤度関数  $\ln L$  は

$$\ln L(\boldsymbol{\beta}) = \sum_{i=1}^N \sum_{j=1}^p y_{ij} \ln \pi_j(\mathbf{x}_i|\boldsymbol{\beta}) + (1 - y_{ij}) \ln (1 - \pi_j(\mathbf{x}_i|\boldsymbol{\beta})) + \boldsymbol{\beta} \text{ を含まない部分} \quad (24)$$

$$\ln L(\boldsymbol{\beta}) = \sum_{i=1}^N \sum_{j=1}^p f_{ij} \ln \pi_j(\mathbf{x}_i|\boldsymbol{\beta}) + (n_{ij} - f_{ij}) \ln (1 - \pi_j(\mathbf{x}_i|\boldsymbol{\beta})) + \boldsymbol{\beta} \text{ を含まない部分} \quad (25)$$

となる。

### 2.2.1 例

計数データ：2項分布

平均（成功確率）： $x_i$  は全て 1

多群の平均値の推定： $x_i$  は群を示すダミー変数（デザイン行列）

ロジスティック回帰：成功確率の対数がパラメタに関して線形

項目反応理論： $X$  は潜在変数（これも推定する）

## 3 ベイズ推定法

パラメタに関する事前の知識がその事前分布  $h(\beta)$  という形で知られている場合には、パラメタの事後分布

$$h(\beta) \propto g(y|X, \beta) \times h(\beta) \quad (26)$$

となるため、その事後モード（同時モード、周辺モード）や事後平均を求めることが出来る。

## 4 最小2乗法

RSS は非負の関数であるから、その最小値が存在し、それを与えるパラメタの値は

$$\frac{\partial \text{RSS}}{\partial \beta} = 0 \quad (27)$$

の根となっているはずである。

### 4.1 線形の場合の例

- 回帰分析
- 冗長性分析
- 主成分分析
- その他

スカラー値を返す行列の関数の微分を用いると便利である。

## 5 交互最小2乗法

複数のパラメタ  $\xi = \xi_1, \xi_2, \dots, \xi_r$  を持つ最小2乗基準をすべてのパラメタに関して同時に最小とするためには

$$\frac{\partial \text{RSS}}{\partial \xi} = 0$$

という同時方程式を解かなければならないが、これは時として困難である。そのような場合には、パラメタを適当に分割し、その分割ごとに RSS の最小化を繰り返していけば、最終的には RSS の極小値が得られることが知られている。すなわち、パラメタの分割を

$$\xi = \{\xi_1, \xi_2, \dots, \xi_s\},$$

その初期値を  $\xi^0$  とすれば,

$s = 1, 2, \dots, S$  に関して,  $\xi_k^0, k \neq s$  を既知とした時に RSS を最小とするような  $\xi_s$  を求め, その値を新たに  $\xi_s^0$  とする.

という過程を繰り返していけば, 最終的に得られる  $\xi^0$  において RSS の極小値が求められる. このような方法は交互最小2乗法 ( Alternating Least Squares method: ALS ) と呼ばれている. ここではその概要を説明する.

RSS は  $\xi$  の関数であるからそれを  $RSS(\xi)$  と書く. また, RSS を  $\xi_s$  のみの関数と見たものを  $RSS(\xi_s)$ ,  $\xi_s$  以外のパラメータを  $\xi_{(s)}$  と書けば, ALS の過程は,  $\xi_{(1)}$  を既知として  $RSS(\xi_1)$  を最小にするような  $\xi_1$  をもとめる, 求められた  $\xi_1$  の値で現在の  $\xi$  を更新する,  $\xi_{(2)}$  を既知として  $RSS(\xi_2)$  を最小にするような  $\xi_2$  をもとめる, 求められた  $\xi_2$  の値で現在の  $\xi$  を更新する,  $\dots$  という過程を繰り返していくことになる.  $RSS(\xi_s)$  は, RSS と基本的には同じものであるから, 各段階は RSS の値そのものを少しずつ改善している. したがって, 各  $\xi_s$  が更新された段階で  $RSS(\xi)$  を評価すれば, 少なくともその値は一つ前の段階で得られた RSS の値よりも大きくはないから, 最終的に, RSS の極小値が得られることになる. すなわち, ALS が収束した段階においては,  $\frac{\partial RSS}{\partial \xi} = 0$  となっている. なお, 収束の判定は, ALS においては各段階における RSS の値が決して大きくならないため, たとえば, RSS の減少分がゼロと見なせるか否かで行えばよい. したがって, ALS のアルゴリズムはおおよそ次のようになる.

```

Store the initial value in  $\xi$ ;
rssp = RSS( $\xi$ );
convergence = 0;
do iteration=1 to maxiter (until convergence = 1);
  do s=1 to S; * for each subset of parameters;
    Get  $\xi_s$  from  $\xi$ ;
    Find the value,  $\hat{\xi}_s$  of  $\xi_s$  which mminizes  $RSS(\xi_s)$ ,
    or, at least, satisfies that  $RSS(\hat{\xi}_s) \leq RSS(\xi_s)$ ;
    Update the corresponding elements of  $\xi$  by  $\hat{\xi}_s$ ;
  end;
  rss = RSS( $\xi$ );
  if rssp - rss  $\leq \epsilon$  then convergence = 1; else convergence = 0;
  rssp = rss;
end;

```

なお, 計算機の容量の増大や速度の改善により, 多くの問題を非線形の最小2乗法 ( section ?? を参照) として直接的にすることが可能となってきたが, そのような場合でも, 条件付き最小2乗解を求めるのが簡単な場合には (多くの場合, 条件付き最小2乗解は線形の問題となる), 交互最小2乗法を用いると効果的な場合が多い. また, 変数の最適変換を含む最小2乗法のように交互最小2乗法を用いざるを得ない問題も多い. 一般に, 非線形の最小2乗法による定式化よりも交互最小2乗法を用いた方が収束までの繰り返し計算の回数は増えるが, 計算時間としてはさほど変わらない場合が多い. また, Ramsay(1975) の方法等を用いれば交互最小2乗法を加速することができる.

以下に,  $v = U\xi + e$  というモデルのパラメータを ALS を用いて推定するプログラムの例を示す. ここでは,  $Q$  個のパラメータのひとつひとつをサブセットとしている. この場合, モデルが,

$$v = U_{(s)}\xi_{(s)} + u_{(s)}\xi_s + \varepsilon$$

ただし,  $u_{(s)}$  は  $U$  行列の第  $s$  列,  $U_{(s)}$  は  $U$  行列の第  $s$  列を除いた  $N \times (Q-1)$  の行列,  $\xi_{(s)}$  は  $\xi$  ベクトルの第  $s$  要素を除いた  $(Q-1) \times 1$  のベクトル, と書けることから, 第  $s$  サブセットのパラメータ  $\xi_s$  の条件付きの最小2乗推定値は,  $(v - U_{(s)}\xi_{(s)})$  を  $u_{(s)}$  に回帰することにより,

$$\hat{\xi}_s = (u_{(s)}'u_{(s)})^{-1}u_{(s)}'(v - U_{(s)}\xi_{(s)})$$

となる.

## 6 非線形モデルの最小2乗解

モデル  $\theta(\xi)$  がパラメタ  $\xi$  ( $Q \times 1$ ) に関して非線形の場合、最小2乗基準

$$\text{RSS}(\xi) = \|\mathbf{v} - \theta(\xi)\|^2$$

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  $r(\xi)$  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  $\xi$  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  $\xi$  の初期値から、 $r(\xi)$  が減少する方向  $\mathbf{d}$  へある分量  $\alpha$  だけ修正する、

$$\xi_{(\ell+1)} = \xi_{(\ell)} + \alpha \mathbf{d}_{(\ell)} \quad (28)$$

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  $r(\xi)$  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

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

に取るものである。この、 $Q \times 1$  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  $\xi_{(\ell)}$  における  $r(\xi)$  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.(28) を用いることにより、

$$r(\xi_{(\ell+1)}) \leq r(\xi_{(\ell)}) \quad (29)$$

を満たす  $\alpha$  が存在することが知られている。この方法は単純であるが、関数の極値に到達するまでの繰り返しの数が多くかかることが欠点である。

一次微分のみならず、2次微分の情報をも用いて最適化を行う方法が、通常 Newton-Raphson 法と呼ばれる方法である。 $r(\xi)$  を  $\xi_{(\ell)}$  のまわりで2次の項まで Taylor 展開すると

$$r(\xi) \approx r(\xi_{(\ell)}) + \mathbf{g}(\xi_{(\ell)})'(\xi - \xi_{(\ell)}) + \frac{1}{2}(\xi - \xi_{(\ell)})' \mathbf{H}(\xi_{(\ell)})(\xi - \xi_{(\ell)})$$

となるが、この  $Q \times Q$  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）

$$\mathbf{H}(\xi_{(\ell)}) = \left[ \frac{\partial^2 r(\xi)}{\partial \xi \partial \xi'} \right]_{\xi=\xi_{(\ell)}}$$

はヘシアン行列と呼ばれるものである。したがって、極値の条件は

$$\frac{\partial r(\xi)}{\partial \xi} \approx \mathbf{g}(\xi_{(\ell)}) + \mathbf{H}(\xi_{(\ell)})(\xi - \xi_{(\ell)}) = 0$$

となり、これより

$$\xi = \xi_{(\ell)} - \mathbf{H}(\xi_{(\ell)})^{-1} \mathbf{g}(\xi_{(\ell)})$$

となる。したがって、Eq.(28) において

$$\mathbf{d}_{(\ell)} = -\mathbf{H}(\xi_{(\ell)})^{-1} \mathbf{g}(\xi_{(\ell)})$$

とおいたものが Newton-Raphson 法の更新式として得られる。Newton-Raphson 法は一般的に効率の良い方法であるが、ヘシアン行列の計算が複雑な場合には利用が困難であること、また、ヘシアン行列が必ずしも非負定値行列とはならないために、まれに Eq.(29) を満たすような  $\alpha$  を見いだすことが不可能な場合が存在する（負のグレイディエント方向から 90 度以上方向がそれる）等の欠点がある。このような欠点を補うために、1 次微分の情報のみからヘシアン行列の逆行列を近似的に計算する方法（いわゆる準ニュートン法）も考案されている。

以下には、最適化されるべき関数が最小 2 乗基準の場合（ $r(\boldsymbol{\xi}) = \text{RSS}(\boldsymbol{\xi})$ ）に適した Gauss-Newton 法に関して説明する。いま、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、

$$\boldsymbol{\theta}(\boldsymbol{\xi}) \approx \boldsymbol{\theta}(\boldsymbol{\xi}_{(\ell)}) + \mathbf{F}(\boldsymbol{\xi}_{(\ell)})(\boldsymbol{\xi} - \boldsymbol{\xi}_{(\ell)})$$

となるが、ここに、

$$\mathbf{F}(\boldsymbol{\xi}) = \frac{\partial \boldsymbol{\theta}(\boldsymbol{\xi})}{\partial \boldsymbol{\xi}} = \left[ \frac{\partial \theta_i(\boldsymbol{\xi})}{\partial \xi_k} \right], \quad i = 1, 2, \dots, N, \quad k = 1, 2, \dots, Q$$

は、ヤコビアン行列と呼ばれる  $N \times Q$  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

$$\mathbf{F}(\boldsymbol{\xi}_{(\ell)}) = \left. \frac{\partial \boldsymbol{\theta}(\boldsymbol{\xi})}{\partial \boldsymbol{\xi}} \right|_{\boldsymbol{\xi}=\boldsymbol{\xi}_{(\ell)}}$$

である。この 1 次近似を RSS に代入すれば、

$$\begin{aligned} \text{RSS}(\boldsymbol{\xi}) &\approx \|\mathbf{v} - (\boldsymbol{\theta}(\boldsymbol{\xi}_{(\ell)}) + \mathbf{F}(\boldsymbol{\xi}_{(\ell)})(\boldsymbol{\xi} - \boldsymbol{\xi}_{(\ell)}))\|^2 \\ &= \|\mathbf{v} - \boldsymbol{\theta}(\boldsymbol{\xi}_{(\ell)}) - \mathbf{F}(\boldsymbol{\xi}_{(\ell)})(\boldsymbol{\xi} - \boldsymbol{\xi}_{(\ell)})\|^2 \end{aligned}$$

となるから、

$$\frac{\partial \text{RSS}(\boldsymbol{\xi})}{\partial \boldsymbol{\xi}} \approx -2\mathbf{F}(\boldsymbol{\xi}_{(\ell)})'(\mathbf{v} - \boldsymbol{\theta}(\boldsymbol{\xi}_{(\ell)})) + 2\mathbf{F}(\boldsymbol{\xi}_{(\ell)})'\mathbf{F}(\boldsymbol{\xi}_{(\ell)})(\boldsymbol{\xi} - \boldsymbol{\xi}_{(\ell)}) = 0$$

の解は

$$\boldsymbol{\xi} = \boldsymbol{\xi}_{(\ell)} + (\mathbf{F}(\boldsymbol{\xi}_{(\ell)})'\mathbf{F}(\boldsymbol{\xi}_{(\ell)}))^{-1}\mathbf{F}(\boldsymbol{\xi}_{(\ell)})'(\mathbf{v} - \boldsymbol{\theta}(\boldsymbol{\xi}_{(\ell)}))$$

となる。また、最適化されるべき関数が最小 2 乗基準の場合には、先に示したグレイディエントベクトルが

$$\mathbf{g} = \frac{\partial \text{RSS}(\boldsymbol{\xi})}{\partial \boldsymbol{\xi}} = -2\mathbf{F}(\boldsymbol{\xi})'(\mathbf{v} - \boldsymbol{\theta}(\boldsymbol{\xi}))$$

となるため、Gauss-Newton 法では Eq.(28) において

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

と置いたものになっていることがわかる。Gauss-Newton 法は 1 次の情報しか利用しないが効率はよく、Eq.(29) 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、

$$\text{RSS} = \|\mathbf{V} - \boldsymbol{\Theta}(\boldsymbol{\xi})\|^2 = \sum_{j=1}^P \|\mathbf{v}_{(j)} - \boldsymbol{\theta}_{(j)}(\boldsymbol{\xi})\|^2$$

と書けるため、変数ごとのヤコビアン行列を

$$\mathbf{F}_j(\boldsymbol{\xi}) = \frac{\partial \boldsymbol{\theta}_{(j)}(\boldsymbol{\xi})}{\partial \boldsymbol{\xi}} = \left[ \frac{\partial \theta_{ij}(\boldsymbol{\xi})}{\partial \xi_k} \right], \quad i = 1, 2, \dots, N, \quad k = 1, 2, \dots, Q$$

と定義すれば、グレイディエント並びに更新の方向は

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

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

で与えられる。

パラメタを更新すべき方向が定まれば、ステップサイズ  $\alpha$  の値を適当に決めることにより、Eq.(29) 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.(29) 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ようになる。

### 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 );

$\mathbf{F}(\boldsymbol{\xi}_{(\ell)})$ ,  $\mathbf{g}(\boldsymbol{\xi}_{(\ell)})$ ,  $\mathbf{d}_{(\ell)}$  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;

## 6.1 非線形の場合の例

- MDS
- ニューラルネットワーク
- その他

## 7 周辺尤度を最大にする方法：EM アルゴリズム

## 8 その他の方法