

最小二乗法に関するメモ*

千葉豪

平成 30 年 1 月 12 日

1 基本的な考え方

1.1 基本中の基本

あるパラメータ x に依存して決まるパラメータ y があるとする¹。 I 個の異なる x の値に対して y の値が得られたものとし、そのセットを (x_1, y_1) 、 (x_2, y_2) 、 \dots 、 (x_I, y_I) と記述することとする。このふたつのパラメータについて $y = ax + b$ なる関係があると仮定し、この I 個のデータセットから係数 a 、 b を決めることを考える。

最小二乗法では、以下で定義される G が最小化されるように、これらの係数を決める。

$$G = \sum_{i=1}^I \{(ax_i + b) - y_i\}^2 \quad (1)$$

G が最小値をとる場合には、以下の式が満足される。

$$\frac{\partial G}{\partial a} = 0, \quad (2)$$

$$\frac{\partial G}{\partial b} = 0 \quad (3)$$

それではこれらの式に式 (1) を実際に代入してみよう。このとき、式 (2) は以下のように書ける。

$$\frac{\partial G}{\partial a} = \sum_{i=1}^I 2(ax_i + b - y_i) \cdot x_i = 2 \left\{ a \sum_{i=1}^I x_i^2 + b \sum_{i=1}^I x_i - \sum_{i=1}^I x_i y_i \right\} = 0 \quad (4)$$

同様に、式 (3) は以下のように書ける。

$$\frac{\partial G}{\partial b} = \sum_{i=1}^I 2(ax_i + b - y_i) = 2 \left\{ a \sum_{i=1}^I x_i + b \sum_{i=1}^I 1 - \sum_{i=1}^I y_i \right\} = 0 \quad (5)$$

式 (4)、(5) は以下のようにまとめられる。

$$a \sum_{i=1}^I x_i^2 + b \sum_{i=1}^I x_i = \sum_{i=1}^I x_i y_i, \quad (6)$$

$$a \sum_{i=1}^I x_i + b \sum_{i=1}^I 1 = \sum_{i=1}^I y_i \quad (7)$$

*/Document/Education/LeastSquare/

¹例としては、GM 計数管の印加電圧 (x に対応) と放射線の計数率 (y に対応) が挙げられるであろう。

2 個の未知数 a 、 b に対して 2 本の方程式が立っているため、これらの未知数を一意的に決めることが出来ることが分かるであろう²。

式 (6)、(7) を行列形式で記述すると以下のように書ける。

$$\begin{pmatrix} \sum_{i=1}^I x_i^2 & \sum_{i=1}^I x_i \\ \sum_{i=1}^I x_i & \sum_{i=1}^I 1 \end{pmatrix} \begin{pmatrix} a \\ b \end{pmatrix} = \begin{pmatrix} \sum_{i=1}^I x_i y_i \\ \sum_{i=1}^I y_i \end{pmatrix} \quad (8)$$

これを行列 M 、ベクトル β 、 n を用いて以下のように記述する。

$$M\beta = n \quad (9)$$

未知数 a 、 b を要素に持つベクトル β は、行列 M に逆行列が存在すれば (行列 M が non-singular であれば) 以下のように一意的に決めることが出来る³。

$$\beta = M^{-1}n \quad (11)$$

もし、 $y = ax + b$ ではなく、 $y = b$ のように y が x に依存しない振る舞いをすると仮定するならばどうなるであろうか。この場合は、式 (5) において $a = 0$ とすればよく、以下の式が得られる。

$$b \sum_{i=1}^I 1 - \sum_{i=1}^I y_i = 0 \quad (12)$$

従って、未知数 b は以下のように決まる。

$$b = \frac{\sum_{i=1}^I y_i}{\sum_{i=1}^I 1} = \frac{\sum_{i=1}^I y_i}{I} \quad (13)$$

これは、 I 個の y_i の平均をとっていることに他ならない。

1.2 基本

これまでの議論に引き続いて、以下のような行列、ベクトルを定義しよう。

$$\mathbf{X} = \begin{pmatrix} x_1 & 1 \\ x_2 & 1 \\ \vdots & \vdots \\ x_I & 1 \end{pmatrix}, \quad \mathbf{y} = \begin{pmatrix} y_1 \\ y_2 \\ \vdots \\ y_I \end{pmatrix} \quad (14)$$

²そうとはならない場合もあるがここでは省略。

³なお、 2×2 の行列の逆行列は以下のように求められる。

$$\begin{pmatrix} a & b \\ c & d \end{pmatrix}^{-1} = \frac{1}{ad - bc} \begin{pmatrix} d & -b \\ -c & a \end{pmatrix} \quad (10)$$

これらの行列、ベクトルを用いて、式 (1) は以下のように書くことができる。

$$G = \sum_{i=1}^I \{(ax_i + b) - y_i\}^2 = (\mathbf{X}\boldsymbol{\beta} - \mathbf{y})^T (\mathbf{X}\boldsymbol{\beta} - \mathbf{y}) \quad (15)$$

実は、式 (9) で定義した M 、 n は \mathbf{X} 、 \mathbf{y} を用いて以下のように記述できる。

$$M = \mathbf{X}^T \mathbf{X} = \begin{pmatrix} x_1 & x_2 & \dots & x_I \\ 1 & 1 & \dots & 1 \end{pmatrix} \begin{pmatrix} x_1 & 1 \\ x_2 & 1 \\ \vdots & \vdots \\ x_I & 1 \end{pmatrix}, \quad (16)$$

$$n = \mathbf{X}^T \mathbf{y} = \begin{pmatrix} x_1 & x_2 & \dots & x_I \\ 1 & 1 & \dots & 1 \end{pmatrix} \begin{pmatrix} y_1 \\ y_2 \\ \vdots \\ y_I \end{pmatrix} \quad (17)$$

従って、最小二乗法で与えられた式 (9) は、 \mathbf{X} 、 \mathbf{y} を用いて以下のように書き換えられる。

$$\mathbf{X}^T \mathbf{X} \boldsymbol{\beta} = \mathbf{X}^T \mathbf{y} \quad (18)$$

ここで、データセット (x_1, y_1) 、 (x_2, y_2) 、 \dots 、 (x_I, y_I) を用いて、以下の方程式を考える。

$$ax_1 + b = y_1, \quad (19)$$

$$ax_2 + b = y_2, \quad (20)$$

$$\dots \quad (21)$$

$$ax_I + b = y_I \quad (22)$$

データセットの個々の値は誤差を含んでばらつくため、この方程式を満足する解は基本的にはあり得ない。さて、これらの式は行列形式では以下のように書ける。

$$\mathbf{X}\boldsymbol{\beta} = \mathbf{y} \quad (23)$$

前掲の式 (18) は、この方程式に対する正規方程式 (Normal equation) と呼ばれる。

2 Meyer の教科書から

Meyer の教科書 [1] における最小二乗法についての記述を以下にまとめる。

観測パラメータ Y が、複数のパラメータ X_n ($n = 1, 2, \dots, N$) の線形結合で表せるとすると、 Y は以下の式で記述される。

$$Y = \beta_1 X_1 + \dots + \beta_N X_N \quad (24)$$

ここで我々は、 X_n と Y の値を観測から得ることにより、未知の定数 (パラメータ) である β_n を求めることを考える⁴。

⁴前節の例 ($y = ax + b$) のように定数項 (b) が存在する場合は、 X_n のうちひとつが 1.0 で固定されていると考えればよいだろう。

さて、 X_n の値が誤差なく求められる（測定される）一方、測定される Y の値は誤差を含んだ形で求められるとする。この場合、式 (24) の代わりに以下の式が書けるであろう。

$$Y = \beta_1 X_1 + \cdots + \beta_N X_N + \epsilon \quad (25)$$

ϵ は測定誤差に対応するランダム変数を示す。従って、 Y はある確率分布に従う確率変数となる。

ここで、 I 回の観測を行ったものとし、 $I > N$ とする。 y_i が i 番目の観測の結果得られた Y の値であるとするならば、 y_i は式 (25) に従い以下のように書くことができるであろう。

$$y_i = \beta_1 x_{1i} + \cdots + \beta_N x_{Ni} + \epsilon_i, \quad i = 1, 2, \dots, I \quad (26)$$

ϵ_i は i 番目の観測における測定誤差である。一般的に、測定誤差は、観測回に依存せずある一定の確率分布に従い、その確率分布の期待値はゼロ、また、観測毎に独立であると見做せる。従って、

$$E[\epsilon_i] = 0, \quad \text{Cov}[\epsilon_i, \epsilon_j] = \begin{cases} \sigma^2 & (i = j), \\ 0 & (i \neq j). \end{cases} \quad (27)$$

と書ける。これらを用いると、観測値 y_i の期待値は

$$E[y_i] = E \left[\sum_{n=1}^N \beta_n x_{ni} \right] + E[\epsilon_i] = \sum_{n=1}^N \beta_n x_{ni} \quad (28)$$

と書ける。

さて、式 (26) は行列形式では $\mathbf{y} = \mathbf{X}_{I \times N} \boldsymbol{\beta} + \boldsymbol{\epsilon}$ と書ける。このとき、式 (28) は $E[\mathbf{y}] = \mathbf{X} \boldsymbol{\beta}$ と書ける。一般的に、測定では、異なるパラメータ X_i と X_j の寄与が全ての測定において同一の比率となるようには行わない⁵。従って、行列 $\mathbf{X}_{I \times N}$ の全ての列は互いに独立であると言え、 $I > N$ を仮定しているので $\text{rank}(\mathbf{X}_{I \times N}) = N$ となる。

今、我々は、観測値 \mathbf{y} からパラメータ $\boldsymbol{\beta}$ を推定することを目的としている。そこで、観測値 \mathbf{y} に対して線形である $\boldsymbol{\beta}$ の予測子 $\bar{\boldsymbol{\beta}} = \mathbf{M}_{I \times M} \mathbf{y}$ を考える。すると、我々が考えるべき問題は、 $\boldsymbol{\beta}$ に対する線形予測子のうち、（分散が最小となる意味で）最良で、不偏な（その期待値が $\boldsymbol{\beta}$ となる）ものが何であるかについて考えることに置き換わる。

まずは、線形予測子が不偏となる条件、すなわち $E[\bar{\boldsymbol{\beta}}] = \boldsymbol{\beta}$ となる条件について考えよう。線形予測子 $\bar{\boldsymbol{\beta}}$ の期待値は

$$E[\bar{\boldsymbol{\beta}}] = E[\mathbf{M} \mathbf{y}] = \mathbf{M} E[\mathbf{y}] = \mathbf{M} \mathbf{X} \boldsymbol{\beta} \quad (29)$$

と書ける。ここで、 $\text{rank}(\mathbf{X}_{I \times N}) = N$ なので、 \mathbf{X} の擬似逆行列⁶ を \mathbf{X}^\dagger とすると $\mathbf{X}^\dagger \mathbf{X} = \mathbf{I}_{N \times N}$ が成り立つ。従って、 $\mathbf{M} = \mathbf{X}^\dagger$ とすれば $E[\bar{\boldsymbol{\beta}}] = \boldsymbol{\beta}$ となることから、線形予測子 $\bar{\boldsymbol{\beta}} = \mathbf{X}^\dagger \mathbf{y}$ が不偏であることが分かる。

次は $\bar{\boldsymbol{\beta}} = \mathbf{X}^\dagger \mathbf{y}$ が最良の線形予測子になることを示そう。不偏の線形予測子として $\boldsymbol{\beta}^* = \mathbf{L} \mathbf{y}$ を仮定すると、前述のように不偏であることから

$$E[\boldsymbol{\beta}^*] = \mathbf{L} E[\mathbf{y}] = \mathbf{L} \mathbf{X} \boldsymbol{\beta} = \boldsymbol{\beta} \quad (32)$$

⁵例えば、全ての測定で $X_j/X_i = \alpha$ であったとするならば、測定データから X_i と X_j の各々に対する β を評価できないことは直感的に理解できるであろう。

⁶行列 $\mathbf{A}_{I \times N}$ の URV 分解を

$$\mathbf{A}_{I \times N} = \mathbf{U} \begin{pmatrix} \mathbf{C} & \mathbf{0} \\ \mathbf{0} & \mathbf{0} \end{pmatrix}_{I \times N} \mathbf{V}^T \quad (30)$$

と書いたとき、擬似逆行列 \mathbf{A}^\dagger は、

$$\mathbf{A}_{N \times I}^\dagger = \mathbf{V} \begin{pmatrix} \mathbf{C}^{-1} & \mathbf{0} \\ \mathbf{0} & \mathbf{0} \end{pmatrix}_{N \times I} \mathbf{U}^T \quad (31)$$

と書ける。ここで、 $\text{rank}(\mathbf{C}) = \text{rank}(\mathbf{A})$ である。

が得られる。 β に依存せずこの式が成り立つためには $LX = I_{N \times N}$ が満足されなければならない。さて、不偏の線形予測子 β^* の n 番目の要素の分散は以下のように計算できる。

$$\text{Var}[\beta_n^*] = \text{Var}[L_{n*} \mathbf{y}] = \text{Var} \left[\sum_{i=1}^I l_{ni} y_i \right] = \sum_{i=1}^I l_{ni}^2 \text{Var}(y_i) = \sigma^2 \sum_{i=1}^I l_{ni}^2 = \sigma^2 \|L_{n*}\|_2^2 \quad (33)$$

$L_{n*} X = e_n^T$ より、 β_n^* の分散を最小にする L_{n*} は、 $z^T X = e_n^T$ を満足する z のうちそのノルムが最小のものであることが分かる。擬似逆行列に関する理論より、それは $z^T = e_n^T X^\dagger = X_{n*}^\dagger$ と与えられることが分かっている。従って、全ての n について β_n^* の分散を最小にするのは $L_{n*} = X_{n*}^\dagger$ となるため、分散を最小にする不偏線形予測子は $\beta = X^\dagger \mathbf{y}$ と与えられることが分かる。

以上をまとめると、以下ようになる。

式 (24) で記述されるように N 個のパラメータ X_n に依存して決まるパラメータ Y があり、これらのパラメータについての I 個の測定データを用いて β を推定する場合 (ただし $\text{rank}(X_{I \times N}) = N$)、不偏かつ分散を最小にするものは以下で与えられる。

$$\beta = X^\dagger \mathbf{y} \quad (34)$$

$\text{rank}(X_{I \times N}) = N$ のとき、

$$X^\dagger = (X^T X)^{-1} X^T \quad (35)$$

なのでこれを式 (34) に代入すると以下を得る。

$$\beta = (X^T X)^{-1} X^T \mathbf{y} \quad (36)$$

これは正規方程式 (18) と対応するものであることが分かるであろう。

参考文献

- [1] C. Meyer, *Matrix analysis and applied linear algebra*, Society for industrial and Applied Mathematics, Philadelphia (2000).