

データ解析のツールとしての チェビシェフの不等式の多次元化の試み

新谷 明雲

山口県立大学 共通教育機構

Multi-Dimensional Chebychev's Inequality as a Tool in Data Analysis

Meiun SHINTANI

The General Education Division of Yamaguchi Prefectural University

Abstract

We provide three types of multi-dimensional Chebychev's inequality corresponding to the choice of domain surrounding the origin in the data space. The first choice of the domain is the hyper-cube, the second one is the hyper-sphere, and the last one is the hyper-ellipse characterized by the correlation coefficient in the L-dimensional data space. The corresponding inequality is respectively given. That each of them has merits and demerits is illustrated

Key words: multi-dimensional Chebychev's inequality, hyper-cube, hyper-sphere, hyper-ellipse, set theory, de Morgan's theorem, correlation coefficient, covariance matrix

キーワード：多次元チェビシェフの不等式、超立方体、超球体、超楕円体、集合論、ド・モルガンの定理、相関係数、共分散行列

§ 1 はじめに

近年、チェビシェフの不等式 [1] は大量のデータを解析する分野においてにきわめて有効なツールとなっていることはつとに知られている [2]。

本稿のねらいは、多変量データを取り扱うために1次元データ列に対するチェビシェフの不等式を多次元化することにある。

まず「 n 個の2次元データ列とは何か」ということから説明しよう。1つの対象に属する2組の変量（例えばある人の身長 x と体重 y ）に注目し、その値を2次元ユークリッド空間の1点に対応させ、これを直交座標（デカルト座標） (x, y) で記す。対象が n 個（ n 人）有れば2次元ユークリッド空間の n 個の点として表示できる。これが n 個の2次元データ列の意味である。本稿では便宜のため生のデータ列 (x, y) よりも無次元化した標準化変数（標準化変数ともいう） (z_x, z_y) を扱うこととする。 x の標準化変数 z_x は $(x - \bar{x})/s(x)$ で定義される。ここで \bar{x} は平均値、 $s(x)$ は x の標準

偏差であり、ともにデータ列（統計量ともいう） x から算出される。

「2次元チェビシェフの不等式」とは、2次元空間内の n 個のデータ点のうち原点を含む決められた閉領域 D の中にどれだけの割合（相対度数）で存在するかを不等式の形で評価することを可能ならしめるものである。同じく L 次元データ列にこれを拡張したものが「 L 次元チェビシェフの不等式」となる。ここで「原点を含む決められた閉領域 D 」とは、いったい何を指すのかあいまいかもしれない。確かにあいまいである。チェビシェフの不等式はこの閉領域 D に応じた数だけ不等式は存在しうるわけである。あなたが好む領域を選びなさい。が好きなように D を選んだからと言って、不等式の導出が可能であるかはわからない。またそれが可能としても直観的に領域の図形を想起できるか、描画し易いか、意味のある結論が引き出せるか、などといった実用面からの不等式の真価が問われよう。

本稿では、閉領域 D として L 次元立方体 ($L = 1$ は

線分、 $L = 2$ は正方形)をまず考え、それに内接する L 次元球体、 L 次元楕円体を扱うこととする。不等式の下限値が同じであれば体積(2次元では面積)が最も小さくなるような領域選択が最適と言えるだろう。

第2節ではチェビシェフの不等式の高次元への拡張を試みる。領域 D の選び方として、最初に正方形の拡張である超立方体を扱う。次に、円の拡張である超球体を、最後に共分散行列から導かれる超楕円体を扱う。それぞれの場合に、2次元データからはじめ3次元、4次元、…、 L 次元へと議論を進める。

第3節は、今後の展望と議論にあてられる。

§2 チェビシェフの不等式の多次元への拡張

2.1 1次元チェビシェフの不等式(復習)

チェビシェフの不等式の多次元化への拡張作業は基本的に1次元のチェビシェフの不等式が土台となるのでその導出過程を復習を兼ね概観してみよう。

次の n 個の1次元(標準化)データ列

$$Z_1, Z_2, \dots, Z_n, \left(Z_i \equiv \frac{x_i - \bar{x}}{s(x)} \right) \quad (1)$$

を考える。 z の平均値 および分散 $s^2(z)$ は

$$\bar{z} = 0, \quad s(z)^2 \equiv \frac{1}{n} \sum_{i=1}^n z_i^2 = 1. \quad (2)$$

今、1次元 z 空間を、 $|z| \leq k$, および $|z| > k$, ($k > 1$)に分ける。分散1の式から、

$$1 = \frac{1}{n} \sum_{i=1}^n z_i^2 = \frac{1}{n} \sum_{|z| > k} z_i^2 + \frac{1}{n} \sum_{|z| \leq k} z_i^2 \geq \frac{1}{n} \sum_{|z| > k} z_i^2 \geq \frac{k^2}{n} \sum_{|z| > k} 1$$

が成り立つので

$$\frac{1}{n} \sum_{|z| > k} 1 \leq \frac{1}{k^2} \quad (3)$$

ないしは、

$$\frac{1}{n} \sum_{|z| \leq k} 1 \geq 1 - \frac{1}{k^2} \quad (4)$$

をうる。(3)、(4)の不等式が成り立つためには

$$k > 1 \quad (5)$$

であることが必要とされる。これが1次元のチェビシェフの不等式である。

冒頭で述べたように L 次元へのチェビシェフの不等式の拡張といっても、 L 次元空間の原点を含む閉領域 D の選び方は無数にありうるがその中でも最も単純な以下の3つの場合、立方体(§2.2)、球体(§2.3)、楕円体(§2.4)、に限定しそれぞれについて不等式の一般化をおこなう。

2.2 1辺の長さが $2\sqrt{Q}$ の L 次元立方体の場合

($L = 2$ の場合)

今、全データの集合を U とし、その部分集合として

$$A = \{(z_{x_i}, z_{y_i}) \mid |z_x| \leq \sqrt{Q}, |z_y| \leq \infty\} \subseteq U, \quad (6)$$

$$B = \{(z_{x_i}, z_{y_i}) \mid |z_x| \leq \infty, |z_y| \leq \sqrt{Q}\} \subseteq U \quad (7)$$

を導入する。ここで正の数 \sqrt{Q} は1より大であるが、相対割合に対する不等式が意味を持つように改めて定義される。閉領域 D は A と B の共通集合、すなわち $D = A \cap B$ として定義する。 D は1辺が $2\sqrt{Q}$ の正方形に入るデータの集合を表し、1次元チェビシェフの不等式から U の真部分集合($D \subset U$)となることが分かる。集合 S の要素の数を $n(S)$ とあらわせば、 $n(U) = n$ となる。各部分集合の要素数の間には次の関係が成り立つ。

$$n(U) = n(A) + n(B) - n(A \cap B) \quad (8)$$

また、

$$(A \cup B) \cup (\overline{A \cup B}) = U, \quad (A \cup B) \cap (\overline{A \cup B}) = \phi$$

であるので、

$$n(U) = n = n(A \cup B) + n(\overline{A \cup B}) \quad (9)$$

が成り立つ。ここで \bar{A} は A の補集合を表す。(8)と(9)より、

$$n(\overline{A \cup B}) = n - n(A) - n(B) + n(A \cap B) \quad (10)$$

を得る。この式より $n(A \cap B)$ の下限値は、 $n(\overline{A \cup B})$ の下限値により達成される。

$$n(\overline{A \cup B}) \geq 0 \quad (11)$$

が、成り立つので、

$$n(A \cap B) \geq -n + n(A) + n(B) \quad (12)$$

をうる。一方、1次元チェビシェフの不等式により、

$$\frac{n(A)}{n} \geq 1 - \frac{1}{\sqrt{Q}}, \quad \frac{n(B)}{n} \geq 1 - \frac{1}{\sqrt{Q}}, \quad \text{for } Q > 1 \quad (13)$$

が成り立つので、(12)式は

$$\frac{n(D)}{n} \geq 1 - \frac{2}{\sqrt{Q}} \quad (14)$$

を得る。これが正方形領域 D に対する2次元チェビシェフの不等式にほかならない。ただし、(14)が意味を持つためには、正方形の一辺の半分(\sqrt{Q})が

$$Q > 4 \quad (15)$$

とならねばならない。

($L = 3$ の場合)

次に一辺の長さが $2k$ の立方体の場合を考えよう。全集合 U の部分集合 A 、 B 、 C を定義する。

$$A = \{(z_{x_i}, z_{y_i}, z_{z_i}) \mid |z_x| \leq \sqrt{Q}, |z_y| \leq \infty, |z_z| \leq \infty\} \subseteq U, \quad (16)$$

$$B = \{(z_{x_i}, z_{y_i}, z_{z_i}) \mid |z_x| \leq \infty, |z_y| \leq \sqrt{Q}, |z_z| \leq \infty\} \subseteq U \quad (17)$$

$$C = \{(z_{x_i}, z_{y_i}, z_{z_i}) \mid |z_x| \leq \infty, |z_y| \leq \infty, |z_z| \leq \sqrt{Q}\} \subseteq U \quad (18)$$

閉領域 D は、共通集合 $D = A \cap B \cap C$ で定義される一辺が $2k$ の立方体である。要素間の関係式として次の関係が成立する。

$$n(A \cup B \cup C) = n(A) + n(B) + n(C) - n(A \cap B) - n(B \cap C) - n(C \cap A) + n(A \cap B \cap C) \quad (19)$$

また、1次元の場合の(4)式と2次元の場合の(14)式を用いることにより以下の式が成立する。

$$\frac{n(A \cap B \cap C)}{n} + \frac{n(\overline{A \cup B \cup C})}{n} \geq 1 - \frac{3}{\sqrt{Q}} \quad (20)$$

この式より $n(A \cap B \cap C)$ の下限値は $n(\overline{A \cup B \cup C})$ の上限値によって達成される。したがって、 $n(\overline{A \cup B \cup C})$ の上限値を求めよう。ド・モルガンの定理により $\overline{A \cup B \cup C} = \overline{A} \cap \overline{B} \cap \overline{C}$ であるので包含関係

$$\overline{A \cup B \cup C} \subset \overline{A}, \overline{B}, \overline{C} \quad (21)$$

が成立するため要素数間に不等式

$$n(\overline{A \cup B \cup C}) \leq n(\overline{A}), n(\overline{B}), n(\overline{C}) \quad (22)$$

が成り立つ。したがって、1次元の1次元チェビシェフの不等式により

$$\frac{n(\overline{A \cup B \cup C})}{n} \leq \frac{1}{\sqrt{Q}} \quad (23)$$

をうる。この式を(20式)に代入して

$$\frac{n(D)}{n} \geq 1 - \frac{4}{\sqrt{Q}} \quad (24)$$

をうる。ここで

$$Q > 16 \quad (25)$$

である。これが立方体D領域に対する3次元チェビシェフの不等式に他ならない。

一般L次元立方体に対しても集合論の定理と一次元チェビシェフの不等式を繰り返し用いることによりL次元チェビシェフの不等式を求めることが可能であると予想される。この点についての詳細報告は稿を改めて紹介することとする。

2. 3 L次元立方体に内接する半径 \sqrt{Q} の超球体の場合

(L = 2の場合)

n個の2次元データ列 (z_{x_i}, z_{y_i}) ($i = 1, 2, \dots, n$) は、各成分が独立に(2)式を満足する。

$$\bar{z}_x = 0, \bar{z}_y = 0 \quad (26)$$

$$s^2(z_x) = \frac{1}{n} \sum_{i=1}^n z_{x_i}^2 = 1, \quad s^2(z_y) = \frac{1}{n} \sum_{i=1}^n z_{y_i}^2 = 1 \quad (27)$$

いま、2次元極座標 (r, θ) を導入し、

$$z_{x_i} = r_i \cos \theta_i, \quad z_{y_i} = r_i \sin \theta_i, \quad (0 < r_i < \infty, 0 \leq \theta_i < 2\pi) \quad (28)$$

(27)式を書き換えると

$$1 = \frac{1}{n} \sum_{i=1}^n r_i^2 \cos^2 \theta_i, \quad 1 = \frac{1}{n} \sum_{i=1}^n r_i^2 \sin^2 \theta_i, \quad (29)$$

となり、この2つの式の辺々を足し合わせると

$$2 = \frac{1}{n} \sum_{i=1}^n r_i^2 (= \bar{r}^2) \quad (30)$$

をうる。この式はrの二乗平均値が2であることを意味する。この値はL次元ではLとなる。この等式を、半径 \sqrt{Q} の内側と外側に分けることを考える。

$$\text{すなわち, } \frac{1}{n} \sum_{i=1}^n r_i^2 = \frac{1}{n} \sum_{r > \sqrt{Q}} r_i^2 + \frac{1}{n} \sum_{r \leq \sqrt{Q}} r_i^2 \geq \frac{1}{n} \sum_{r > \sqrt{Q}} r_i^2 \geq \frac{Q}{n} \sum_{r > \sqrt{Q}} 1$$

が成り立つので

$$\frac{1}{n} \sum_{r > \sqrt{Q}} 1 \leq \frac{2}{Q} \quad (31)$$

ないしは、

$$\frac{1}{n} \sum_{r \leq \sqrt{Q}} 1 \geq 1 - \frac{2}{Q} \quad (32)$$

をうる。ただし

$$Q > 2 \quad (33)$$

を満たす必要がある。これが2次元球面(円)を閉領域Dとするチェビシェフの不等式に他ならない。

(32)式の右辺は、正方形に対するチェビシェフの不等式のそれに等しい。したがって、正方形とそれに内接する円が、同じ評価式を与えるので、面積の小さな円の方が評価式としてはより強力と言える。

(L ≥ 3の場合)

(32)式のL次元超球体に対する拡張は容易で、

$$\frac{1}{n} \sum_{r \leq \sqrt{Q}} 1 \geq 1 - \frac{L}{Q} \quad (34)$$

として与えられる。ただしQに対する制限は(33)式を

$$Q > L \quad (35)$$

で置き換えればよい。3次元においては、立方体の評価式は(24)式で与えられるので、球体の方がより大きな下限値を与える。したがって、球体の方が不等式としては強力と言える。4次元以上でも同様の結論となることが予想されるが、検証は今後の研究に委ねることとしよう。

2. 4 1辺 $2\sqrt{Q}$ のL次元立方体に内接する超楕円体の場合

これまでは閉領域Dの選び方は変量間の相関の強さとは無関係な選択であった。この節では相関の強さを考慮に入れた領域Dに対するチェビシェフの不等式を導こう。

(L = 2の場合)

n個の2次元データ列 (z_{x_i}, z_{y_i}) ($i = 1, 2, \dots, n$) に対し(26)、(27)式がなりたつ。xとyの間の相関の強さを表す統計量として相関係数 r_{xy} を導入する。

$$r_{xy} = r(x, y) = r(z_x, z_y) = \frac{1}{n} \sum_{i=1}^n z_{x_i} z_{y_i}, \quad (-1 \leq r_{xy} \leq 1). \quad (36)$$

いま、次の変量qを考えよう。

$$q_i \equiv \frac{1}{1-r_{xy}^2} \{z_{x_i}^2 - 2r_{xy} z_{x_i} z_{y_i} + z_{y_i}^2\}, \quad (i = 1, 2, \dots, n). \quad (37)$$

この量は、判別式が負となることから z_{x_i}, z_{y_i} の値によらず常に正値を取る。平均値 \bar{q} は

$$\bar{q} = \frac{1}{n} \sum_{i=1}^n q_i = 2 \quad (38)$$

で、円の場合の(30)式に酷似する。とくに無相関のとき、すなわち $r_{xy} = 0$ の場合

$$q_i = r_i^2 \quad (39)$$

となり円の半径の2乗値に一致する。いまこの2次元空間に次の楕円の方程式を考える。

$$Q = \frac{1}{1-r_{xy}^2} \{z_x^2 - 2r_{xy} z_x z_y + z_y^2\}, \quad Q > 0 \quad (40)$$

Q値は常に正であり楕円の大きさを表す量である。式(38)を、Q値より大きい領域と小さい領域に分けて考えよう。

$$2 = \frac{1}{n} \sum_{i=1}^n q_i = \frac{1}{n} \sum_{q > Q} q_i + \frac{1}{n} \sum_{q \leq Q} q_i \geq \frac{1}{n} \sum_{q > Q} q_i \geq \frac{Q}{n} \sum_{q > Q} 1 \quad (41)$$

となるので、次の不等式をうる

$$\frac{1}{n} \sum_{q>Q} 1 \leq \frac{2}{Q} \quad (42)$$

ないしは、

$$\frac{1}{n} \sum_{q \leq Q} 1 \leq 1 - \frac{2}{Q} \quad (43)$$

ただし、Qは

$$Q > 2 \quad (44)$$

を満たさなければならない。これが楕円領域に対する2次元チェビシェフの定理に他ならない。(40)式であたえられる楕円が1辺 $2\sqrt{Q}$ の正方形に内接する楕円であることは、座標軸を 45° 回転することにより理解できる。

$$z_x = \frac{\sqrt{2}}{2}(z'_x + z'_y),$$

$$z_y = \frac{\sqrt{2}}{2}(z'_x - z'_y)$$

として座標変換すると、(40)式は楕円の標準形

$$\frac{z_x'^2}{1+r_{xy}} + \frac{z_y'^2}{1-r_{xy}} = Q \quad (45)$$

となる。Fig. 1にみられるように正の相関 ($r_{xy} > 0$) の場合には長軸は $\sqrt{(1+r_{xy})Q}$ で、短軸が $\sqrt{(1-r_{xy})Q}$ となる。原点と正方形と内接する点を結べば楕円長軸を挟んだ2本の直線

$$z_y = r_{xy}z_x, \quad z_y = \left(\frac{1}{r_{xy}}\right)z_x$$

が得られるが、傾きが r_{xy} のものを「 z_x の z_y への回帰直線」、傾きが $1/r_{xy}$ のものを「 z_y の z_x への回帰直線」とそれぞれ呼ばれる。楕円の面積は

$$\pi Q \sqrt{1-r_{xy}^2} \quad (46)$$

となり、一方1辺が $2\sqrt{Q}$ で面積が $4Q$ の正方形に内接する円の面積は

$$\pi Q (\approx 3.14 Q) \quad (47)$$

であるので、同じ下限値に対し面積は楕円の方が $\sqrt{1-r_{xy}^2}$ だけ小さくなる。とくに、 r_{xy} の絶対値が1に近づくと面積は急激に減少する。したがって、等面積であるならば、楕円領域のチェビシェフの不等式はより大きな下限値を与え、より強力な評価式を我々に提供することになる。

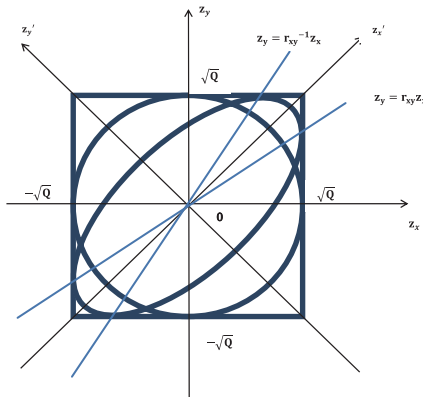


Fig. 1 1辺が $2\sqrt{Q}$ の正方形に内接する円と楕円($r_{xy} > 0$ の場合)

($L \geq 3$ の場合)

2次元の場合と同様に、行うことが可能である。Q値の定義式は多次元正規分布の指数eの肩が $(-\frac{Q}{2})$ であたえられる。すなわちQは共分散行列Vとデータを表す標準化された列ベクトルzから作られる2次形式に他ならない。

$$Q = \mathbf{z}^T \mathbf{V}^{-1} \mathbf{z} \quad (48)$$

ここで、 \mathbf{V}^{-1} は、共分散行列Vの逆行列を表し、zを列ベクトル表記としたら \mathbf{z}^T は行ベクトルに対応する。ちなみに3次元の共分散行列Vは以下のように入力される。

$$\mathbf{V} = \begin{pmatrix} 1 & r_{xy} & r_{xz} \\ r_{xy} & 1 & r_{yz} \\ r_{xz} & r_{yz} & 1 \end{pmatrix} \quad (49)$$

L次元超楕円体に対するチェビシェフの不等式は、

$$\frac{1}{n} \sum_{q \leq Q} 1 \leq 1 - \frac{L}{Q} \quad (50)$$

となる。ただし、

$$Q > L \quad (51)$$

であたえられる。

§ 3 おわりに

前節で、多次元チェビシェフの不等式を、超立方体、超球体、超楕円体について求めることができた。それぞれに長短がある。相関係数を考慮に入れたチェビシェフの不等式は、大量の多変量データを数値解析する研究者向けには最適と言える。特に相関係数の絶対値が1に近ければ近いほどその威力は絶大となる。一方、領域内の数のカウントという意味では、立方体が最も手軽であり、市販の表計算ソフトなどで簡単にカウントが可能である。球体はその中間に位置し、高度のプログラムは必要としないがそれなりの覚悟が必要となる。

L次元の直方体に対する不等式の一般化もこの論文に基づき容易に行うことができる。読者への練習問題として残しておく。

4次元以上の超立方体の場合のチェビシェフの不等式については、別の稿に残すこととする。

参考文献

- [1] 例えば、西田敬義ほか著、大学演習「数理統計学(第8版)」、裳華房(昭和49年)
- [2] 北原知就ほか、「一般化チェビシェフ不等式とその最適化への応用」、数理解析研究所考究録、第1534巻(2008)21-24。この論文の参考文献にチェビシェフの不等式の一般化に関する文献が見られる。