

多項分布モデルによる日降水量頻度分布の同一性の検討とグループ分け

Checking Homogeneity of Frequency-distributions of Daily Precipitation

and Their Practical Grouping by Multinomial Distribution Model and AIC

鈴木正人*・長尾正志**

By Masato SUZUKI and Masashi NAGAO

Steadiness and homogeneity of hydrological quantities are mainly assumed in a customary stochastic analysis. This study gives reconsideration to this assumption, through checking homogeneity of shape of frequency-distributions. Judgement of homogeneity is given by AIC on the basis of multinomial-distribution model. If every distributions are judged not homogenous, practical grouping of those distributions is proposed by the same method.

Numerical calculation is carried out with daily precipitation in Gifu city from 1891 to 1990. It is shown that six-parts dividing model in order of maximum daily precipitation is most suitable.

Keywords:frequency-distribution, homogeneity, multinomial-distribution, AIC

1. はじめに

従来、水文量を統計的に取り扱う際には、おもに資料の定常性・均質性を前提とし、資料の母集団は一つと考えて、確率水文量などの計算が行われてきた。もちろん異常値検定に代表される異質な資料を除外する試みはなされてきたが、あくまでも、主体となる母集団は唯一であるという視点に立脚したものといえよう。そこで、本研究は、日降水量の頻度分布を対象として、従来”前提とする”の一言のもとに片づけられていた水文資料の均質性・定常性の仮定を見直し、均質でないと判断された場合は、実用的に同一とみなせるもの同士にグループ分けすることで、均質性を満たすような水文資料の分割を試みたものである。

2. 多項分布モデルによる分布の同一性の判定

複数の頻度分布が与えられた場合、多項分布モデルを基礎にそれらの頻度分布の同一性をA I Cにより判定する手法が提案されている¹⁾。その詳細は文献を参考にしていただくとして、ここでは周辺・同時の両分布についてその要点を記述しておく。

* 正会員 工博 岐阜工業高等専門学校助手 土木工学科
(〒501-04 岐阜県本巣郡真正町)

** 正会員 工博 名古屋工業大学教授 工学部社会開発工学科
(〒466 名古屋市昭和区御器所町)

2. 1 周辺分布の同一性の判定手法²⁾

n 回の試行が独立で、各回の試行の結果が c 個の排反な事象に分けられる場合、1番目の事象 E_1 が k_1 回、2番目の事象 E_2 が k_2 回、……、 c 番目の事象が k_c 回起る確率は次式の多項分布で与えられる。

$$m(k_1, \dots, k_c | p_1, \dots, p_c) = \frac{n!}{k_1! \dots k_c!} p_1^{k_1} \dots p_c^{k_c} \quad (1)$$

パラメータ p_1, \dots, p_c はそれぞれ対応する事象の生起確率であり、その最尤解は次式で与えられる。

$$p_i = n_i/n \quad (i=0, 1, \dots, c), \quad n_i: \text{事象 } i \text{ の生起度数}, \quad n = \sum n_i \quad (2)$$

以下の記号を用い、各年の水文量頻度分布を多項分布でモデル化することで分布の同一性の検討を行う。

n : 全資料数（計算対象期間の全日数）

$n(i_1, i_2)$: 各年毎の水文量の頻度分布 (i_1 : 年, i_2 : 頻度分布の階級, $i_1=1, 2 \dots, c_1$; $i_2=1, 2, \dots, c_2$)

$n(i_1)$: i_1 年の資料数 (365 または 366)

$p(i_2 | i_1)$: i_1 年に i_2 階級の水文量が生起する確率 ($\sum_{i_2} p(i_2 | i_1) = 1$)

いま、ある資料 $\{n(i_1, i_2)\}$ が得られる確率 ($= P[\{n(i_1, i_2)\} | \{p(i_2 | i_1)\}]$) は、次式のように c_2 個の項を持つ c_1 個の多項分布の積で与えられる。

$$P[\{n(i_1, i_2)\} | \{p(i_2 | i_1)\}] = \prod_{i_1=1}^{c_1} \left\{ \frac{n(i_1)!}{\prod_{i_2=1}^{c_2} n(i_1, i_2)!} \prod_{i_2=1}^{c_2} p(i_2 | i_1)^{n(i_1, i_2)} \right\} \quad (3)$$

ここで、 $p(i_2 | i_1)$ をパラメータとみなしたときの対数尤度は次式で表現される。

$$\begin{aligned} L[p(i_2 | i_1)] &= K_4 + \sum_{i_1} \sum_{i_2} n(i_1, i_2) \ln p(i_2 | i_1) \\ K_4 &= \ln \left\{ \prod_{i_1=1}^{c_1} n(i_1)! / \prod_{i_1=1}^{c_1} \prod_{i_2=1}^{c_2} n(i_1, i_2)! \right\} \end{aligned} \quad (4)$$

以下ではパラメータ $p(i_2 | i_1)$ の与えかたを変えて、それぞれ異なるモデルを表現し、分布の同一性の判定を行う。すなわち、各年が同一の分布に従うとするモデル MODEL (1)，各年がそれぞれ異なる分布に従うとするモデル MODEL (c1)，全体が M 個 ($1 \leq M \leq c_1$) の分布に従うとするモデル MODEL (M) のパラメータをそれぞれつぎのように表現する。

MODEL (1) : $p(i_2 | i_1) = \theta(i_2)$ 自由パラメータ数: $c_2 - 1$ 個

MODEL (c1) : $p(i_2 | i_1) = \theta(i_2 | i_1)$ // : $c_1 \times (c_2 - 1)$ 個

MODEL (M) : $p(i_2 | i_1) = \theta(i_2 | m_j)$ ($j=1, 2, \dots, M$) // : $M \times (c_2 - 1)$ 個

$$\theta(i_2 | m_j) = \sum_{i_1 \in m_j} n(i_1, i_2) / \sum_{i_1 \in m_j} n(i_1), \quad m_j: \text{分割されたグループ} \quad (5)$$

上式の MODEL (M) で $M=1, c_1$ とすると、それぞれ MODEL (1), MODEL (c1) となる。

各モデルの AIC は次式で求められる。

$$AIC = -2 \times L[p(i_2 | i_1)] + 2 \times \text{自由パラメータ数} \quad (6)$$

対数尤度 $L[p(i_2 | i_1)]$ は、基本的には (4) 式で求められるが、各モデルとも K_4 の項は共通なので、実際の計算には K_4 を減じた右辺第 2 項を用いた。

各モデルの AIC を計算し、MODEL (1) が採択されれば、分布は同一であると判断される。また、MODEL (M) が採択されれば、分布は同一でないと判断されるとともに、対応した資料分割数、および分割方法が得られる。

2. 2 同時分布の同一性の判定手法³⁾

日降水量を計算対象として用いる場合、その生起量の持続性が無視できない場合がある。そこで一次の自己相関を考慮した同時頻度分布に対しても、その同一性の判定手法を記述しておく。以下の記号を採用する。

$n(i_1, i_2, i_3)$: i_1 年に i_2, i_3 階級の日降水量が継続して生起した頻度

$$(i_1 = 1, 2, \dots, c_1 : i_2, i_3 = 1, 2, \dots, c_2)$$

$p(i_2, i_3 | i_1)$: i_1 年という条件のもとで (i_2, i_3) 階級の日降水量が継続して生起する確率

$$\sum_{i_2} \sum_{i_3} p(i_2, i_3 | i_1) = 1$$

$p(i_1, i_2, i_3)$: i_1 年, (i_2, i_3) 階級の同時生起確率

$$\sum_{i_1} \sum_{i_2} \sum_{i_3} p(i_2, i_3, i_1) = 1$$

いま, ある資料 $\{n(i_1, i_2, i_3)\}$ が得られる確率 ($= P(\{n(i_1, i_2, i_3)\} | \{p(i_1, i_2, i_3)\})$) はつきの多項分布で表される。

$$P(\{n(i_1, i_2, i_3)\} | \{p(i_1, i_2, i_3)\}) = \frac{n!}{\prod n(i_1, i_2, i_3)!} \prod p(i_2, i_3, i_1)^{n(i_1, i_2, i_3)} \quad (7)$$

ここで, \prod は i_1, i_2, i_3 についての総乗を意味する。 (7) 式より, $p(i_2, i_3 | i_1)$ をパラメータとみなしたときの対数尤度を, $p(i_2, i_3 | i_1)$ に着目して記述すると, 次式となる。

$$L[p(i_2, i_3 | i_1)] = K_6 + \sum_{i_1} \sum_{i_2} \sum_{i_3} n(i_1, i_2, i_3) \ln p(i_2, i_3 | i_1) \\ K_6 = \ln \left\{ \prod_{i_1=1}^{c_1} n(i_1)! / \prod_{i_1=1}^{c_1} \prod_{i_2=1}^{c_2} \prod_{i_3=1}^{c_3} n(i_1, i_2, i_3)! \right\} - \ln n! - \sum \ln n(i)! + \sum n(i) \ln p(i_1) \quad (8)$$

周辺分布の場合と同様に, パラメータ $p(i_2, i_3 | i_1)$ の与えかたによりモデルの表現を行うが, その際, 同時生起確率を独立と考える, すなわち同時分布を周辺分布の積で与えるか, または従属と考えるかでパラメータの表現は異なる。周辺分布の場合と同様に, 一般的なMODEL (M) について, 両者を記述する。まず, 独立とするモデルのパラメータは,

$$p(i_2, i_3 | i_1) = p(i_2 | i_1) \times p(i_3 | i_1) \\ = \theta(i_2 | m_j) \times \theta(i_3 | m_j) \quad \text{自由パラメータ数: } M \times (c_2 - 1) \text{ 個} \\ \theta(i_2 | m_j) = \sum_{i_1 \in m_j} \sum_{i_3=1}^{c_2} n(i_1, i_2, i_3) / \sum_{i_1 \in m_j} n(i_1) \quad (9)$$

また, 従属とするモデルのパラメータは,

$$p(i_2, i_3 | i_1) = \theta(i_2, i_3 | m_j) \quad \text{自由パラメータ数: } M \times (c_2 \times c_2 - 1) \text{ 個} \\ \theta(i_2, i_3 | m_j) = \sum_{i_1 \in m_j} n(i_1, i_2, i_3) / \sum_{i_1 \in m_j} n(i_1) \quad (10)$$

と表現される。周辺分布の場合と同様に, 対数尤度のうち, 共通する部分である K_6 は除外し, AIC によりモデルを採択する。従属モデルは, 独立モデルに比べて, パラメータ数が階級数の2乗の大きさで増加するので採択されにくくなる。

3. 分布の同一性の計算例

周辺分布モデルを例として, 分布の同一性の判定の計算例を示し, 判定にAICを用いる意義を明確にしておく。表-1に示す3種の頻度分布を用いる。頻度分布Aは, 岐阜市の1986年の日降水量を離散化幅15mmで離散化したもの, また頻度分布Bは同じく1971年のものである。また, 頻度分布Cは頻度分布Aの度数を各階級とも5倍したもので, 相対頻度は頻度分布Aと同一である。

まず, 頻度分布A, Bの同一性を判定する。二つの分布を同一とするMODEL (1)の場合, 対数尤度は次式で求められる。

表-1 計算例に用いる頻度分布(相対頻度)

階級	頻度(相対頻度)		
	A	B	C
1	330 (0.904)	321 (0.880)	1650 (0.904)
2	19 (0.052)	27 (0.074)	95 (0.052)
3	5 (0.014)	7 (0.019)	25 (0.014)
4	6 (0.016)	2 (0.006)	30 (0.016)
5	4 (0.011)	1 (0.003)	20 (0.011)
6	0 (0)	1 (0.003)	0 (0)
7	1 (0.011)	3 (0.008)	5 (0.011)
8	0 (0)	1 (0.003)	0 (0)
9	0 (0)	2 (0.006)	0 (0)

$$LL = (330+321) \times \ln\{(330+321)/730\} + (19+27) \times \ln\{(19+27)/730\} + \cdots + (0+2) \times \ln\{(0+2)/730\} = -357.862$$

自由パラメータ数は階級数から1を減じた8だからこのモデルのAIC, AIC(1) = 731.722 となる。

また、二つの分布が異なるとするMODEL(2)の対数尤度は、次式で求められる。

$$LL = 330 \times \ln\{330/365\} + 19 \times \ln\{19/365\} + \cdots + 321 \times \ln\{321/365\} + \cdots + 2 \times \ln\{2/365\} = -351.627$$

自由パラメータ数は、分布を2つと仮定しているのでMODEL(1)の2倍の16だから、AIC(2) = 735.254となり、AIC(1)の方が小さいことから、MODEL(1)が採択され、分布は同一と判定される。

つぎに、頻度分布B, Cの同一性を検討する。先述と同じ手順により、分布を同一と仮定した場合は、LL = -1002.283で、AIC(1) = 2020.568、分布が異なるとするモデルでは、LL = -975.122でAIC(2) = 1982.244となり、MODEL(2)が採択され、分布は同一でないと判定される。

このように、相対頻度が等しい分布A, Cでも、比較の対象となる分布により同一性の判定は異なる。これは、分布A, Bの比較では、対数尤度はMODEL(2)の方が大きく、適したモデルと考えられるが、自由パラメータ数が2倍になるので、対数尤度の差よりも自由パラメータ数の差の方が大きくなり、AICではMODEL(1)が採択される。一方、分布B, Cの比較では、やはり対数尤度はMODEL(2)の方が大きく、自由パラメータ数もMODEL(2)の方が多いのだが、頻度分布の度数が多いため、対数尤度の差が自由パラメータ数の差よりも大きく、AICではMODEL(2)が採択された。

つまり、頻度の度数が大きくなると、対数尤度の絶対値は増加するが、一方自由パラメータ数は変わらないため、AICにおいて対数尤度が占める割合が増加し、分布形状の相違が同一性の判定に支配的に働く。逆に、分布形状にある程度の相違はあっても、度数が少ないと対数尤度の差は小さくなり、分布は同一と判定されやすくなる。したがって、分布の同一性の判定は、分布形状の相違とその生起度数とを総合的に判断して行われているといえよう。

4. 適用計算

4. 1 計算方法

1891年～1990年の岐阜市日降水量を資料として、各年の日降水量の頻度分布（周辺分布、および同時分布）を同一性の検討対象とした。また、全資料数からステージエスの式により、頻度分布の階級数は18個、階級幅は最大値、最小値より15mmとした。

分布の同一性の判定には、2.で述べたように各年の降水量頻度分布の資料についてグループ数を仮定し、幾つかのグループに分割する必要がある。その方法として、無作為に100年分の資料を分割するとすれば、全体を2組に分割するだけでも、 $(2^{100} - 2) / 2$ 通りの分割の方法があり膨大な計算量になる。そこで、本研究では、計算の便宜を計り、基礎統計量を各年の頻度分布の特性を表す指標として用い、資料を年を一括した単位として昇順に順位づけし、その順位にしたがって2組、および3組に分割することにした。2組に分割する場合は、まず1位～100位までを順位づけしておいて、1位と[2～100位]、[1～2位]と[3～100位]、…、[1～99位]と100位、というように99通りの分割の仕方がある。また、3組に分割する場合は、同様に4851 ($=_{99}C_2$) 通りの分割の仕方がある。2または3分割されたグループをさらに分割することで、8分割モデルまでを検討の対象にした。

各分割モデルにおけるAICを比較すれば、その指標における最適な分割が求まる。また各指標における最適分割モデルのAICを比較することから、どの指標による分割が適当かが判断できる。

4. 3 結果および考察

(a) 周辺分布

順位づけの指標には、2分割には、平均、分散、変動係数、歪係数、無降雨日数、年最大日降水量を、3分割には2分割の結果を参考にして、平均、分散、年最大日降水量を指標として用いた。各指標による最適

分割モデルのAICを表-2に示す。表中で、指標に関わらず1分割がMODEL(1)に、100分割がMODEL(100)に対応している。また、6分割において3-2は3分割を2分割、2-3は2分割を3分割したことを意味する。

	平均	分散	歪係数	最大	無降雨	年代順
1分割	36808.34	同左	同左	同左	同左	同左
2分割	36744.78	36742.25	36752.88	36731.93	36807.88	36813.64
3分割	36728.65	36719.78	36724.37	36677.84	—	—
4分割	36729.16	36713.78	36724.92	36652.72	36831.09	—
6分割	3-2 2-3	36759.00 36754.33	36737.83 36733.12	36745.25 36741.23	36656.58 36642.67	— —
8分割	36787.95	36770.18	36781.85	36654.42	—	—
100分割	39155.62	同左	同左	同左	同左	同左

*6分割で3-2は、3分割をそれぞれ2分割、2-3は2分割を3分割したもの

AICが最小なのは、年最大

表-3 最大値を指標とした最適分割の相対頻度分布(頻度分布)

日降水量を指標とした

i	up	middle1	middle2	middle3	middle4	low
1	0.88(4161)	0.89(4541)	0.89(5194)	0.89(4193)	0.89(5840)	0.89(8497)
2	0.07(312)	0.05(277)	0.06(354)	0.06(305)	0.07(433)	0.06(537)
3	0.02(116)	0.03(145)	0.03(147)	0.03(131)	0.02(157)	0.03(256)
4	0.02(80)	0.01(61)	0.01(60)	0.01(58)	0.01(69)	0.01(115)
5	0.01(31)	0.01(39)	0.01(32)	0.00(23)	0.01(36)	0.00(60)
6	0.00(10)	0.00(15)	0.00(20)	0.00(18)	0.00(16)	0.00(34)
7	0.00(10)	0.00(12)	0.00(14)	0.00(4)	0.00(22)	0
8	0.00(7)	0.00(1)	0.00(7)	0.00(15)	0	0
9	0.00(1)	0.00(5)	0.00(18)	0	0	0
10	0	0.00(10)	0	0	0	0
11	0	0.00(5)	0	0	0	0
12	0.00(6)	0	0	0	0	0
13	0.00(3)	0	0	0	0	0
14	0.00(5)	0	0	0	0	0
15	0.00(2)	0	0	0	0	0
16	0.00(1)	0	0	0	0	0
17	0.00(2)	0	0	0	0	0
18	0.00(2)	0	0	0	0	0

年代順に並べた資料を2分割した場合は、MODEL(1)のAICより値が大きいことから、経年にみて日降水量の頻度分布形状にトレンドのような一方的な変化は認められない。

つぎに、年最大日降水量を指標とした場合の8分割の手順を図-1に示す。まず全体を2分割する場合、1~57位、58~100位の分割がAIC(2)=36731.93で最も適当な分割となる。このAIC(2)をMODEL(1)のAIC(1)(=36808.34)と比較すると、AIC(2)の方が小さいので2分割モデルの方が適当と判断され、つぎの4分割へと計算を進める。以後、AIC(1)、AIC(2)をそれぞれ比較し、分割するか否かの判定をしていくことで、この手順における最も適した分割が得られる。図中で太線で囲まれた順位が最終的に得られた分割の結果(1~26、27~44、45~57、58~73、74~87、88~100位の6分割)で、4分割から8分割へのモデルでAICの値の小さい方(図中で*を付したもの)を総計したものが、この分割モデルのAICになる。

ちなみに、この手順で得られた分割方法、AICは、表-2に示した年最大日降水量を指標とした6分割(2-3分割)の結果と一致している。分割された各グループの相対頻度分布および頻度分布を表-3に示す。最大階級は段階的に減少しているが他の階級の相対頻度には明確な変化は認めら

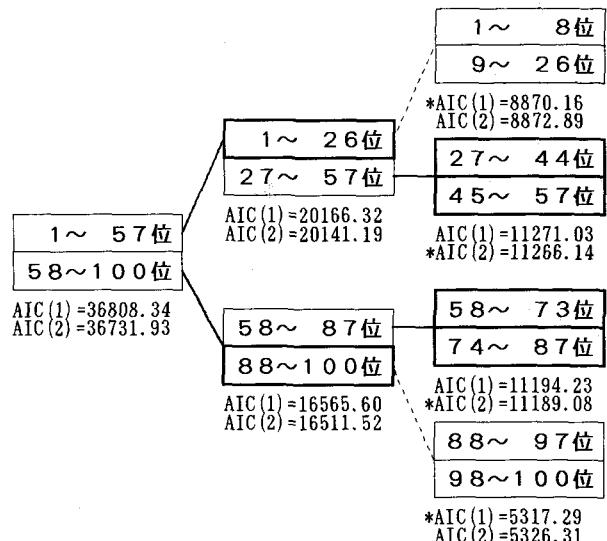


図-1 最大値を指標とした分割の手順(8分割)

れないので、本研究で用いたモデルでは頻度分布の階級の最大値、すなわち年最大日降水量が頻度分布の同一性の判定に支配的に関与していると思われる。

(b) 同時分布

各年日降水量の同時頻度分布に対し、周辺分布への適用で有効だった年最大日降水量と、自己相関係数の2つの指標で順位づけし、2分割を繰り返すことで8分割モデルまでの検討を行った。なお、2.2で述べたように、同時分布の場合は、自己相関性を考慮するか否かで2通りのモデルが作れるので、2分割モデルの場合、分割と自己相関の有無の組み合わせで、計4通りのモデルが検討対象となる。

自己相関係数を指標とした場合は、1~20位（独立）、21~29位（独立）、30~100位（従属）の3分割がAIC=73865.22で最適なモデルという結果が出た。一方、年最大日降水量を指標とした場合は1~8位、9~25位、26~44位、45~56位、57~73位、74~87位、88~92位、93~100位（いずれも独立）の8分割がAIC=73509.76で最適なモデルとなり、自己相関係数を指標として用いた場合よりもAICが小さくなつた。自己相関性の有無は、1~100位を一まとめに取り扱つたモデルでは、従属モデルが採用されたが、他の場合はいずれも独立モデルが採用された。これは、資料を分割することで、頻度分布の度数が減少し、AICのうち自由パラメータ数の増加が対数尤度の増加を上回つてしまい、従属モデルが採択されなかつたのであろう。なお、この分割結果を周辺分布の分割結果である図-1と比較すると、最終的な分割数に差異はあるものの分割の境界の順位はほぼ類似していることがわかる。

(c) 多雨期降水量の周辺分布

以上の各年日降水量の頻度分布に対する検討では、周辺・同時分布とともに年最大日降水量を指標とした分割が適していることがわかつた。そこで、年間の頻度分布ではなく、年最大日降水量が生起した期間（便宜上多雨期と呼ぶ）の周辺分布に対して、同様の適用を行つた。期間の決め方としては、100年の日降水量資料に対し、最大値が生起した月日を調べ、その最早日、最遅日ではさまれた期間を多雨期とした。具体的には、2月16日（1922年に生起）~10月21日（1980年に生起）の237日間である。順位づけの指標には、年最大日降水量、分散、平均の三種を用いたが、最終的に得られた最適モデルは、年最大日降水量を指標とした場合の、1~25位、26~44位、45~72位、73~87位、88~100位の5分割で、AIC=29982.87であった。この結果を各年単位の周辺分布に対する結果と比べると、多雨期の45~72位に相当するグループが、年単位では45~57位、58~72位に2分割されている他は完全に一致した結果となつてゐる。これは、今回用いたモデルでは、年単位の頻度分布も、多雨期の頻度分布も、共にその同一性の判定は年最大日降水量という最大階級により左右されているからであろう。

5. まとめ

頻度分布の同一性の検討という観点から、本研究では、一地点の日降水量頻度分布に対して検討を行つた。その結果、周辺分布、同時分布、ともに年最大日降水量を指標にして順位づけした資料の分割が有効であることがわかつた。多項分布モデルという比較的パラメータ数の多い分布をモデルの基礎として用いたので、同時分布では自己相関性を考慮したモデルは採択されにくかつたが、二項分布、指数分布などのパラメータ数の少ない分布を用いることで、自己相関性の表現も可能になるであろう。今後は、本手法を異なる地点間の頻度分布の同一性の判定にも適用・検討していく予定である。

6. 参考文献

- 1) 坂元慶行・石黒真木夫・北川源四郎：情報量統計学，共立出版，1983,pp.74~77
- 2) 鈴木正人・長尾正志：AICによる水文量頻度分布の同一性の検討，第47回土木学会年次学術講演会概要集，1992,第II部門，pp.670~671
- 3) 前掲1），pp.96~105