

離散分布の合成による渴水時流況のモデル化の研究

Stochastic modeling by the combination of two discrete distributions
of flow-series in dry season

名古屋工業大学博士後期課程学生 鈴木正人 Masato SUZUKI
名古屋工業大学工学部 長尾正志 Masashi NAGAO

Flow series is considered by combination of the following two components. The one for the domain of smaller quantity of flow is a simultaneous correlated binomial distribution, the other for the larger quantity of flow is an independent binomial or Poisson distribution. The whole flow model is composed of the linear combination of the distribution functions for those two models. The decision of boundary for two components and the selection of optimal model are judged by AIC(Akaike's Infomation Criteria). Through the analysis on the simulated data, the estimated results for the model are well coincident with the original parameters. In addition, by applying the model for the observed flow data at the drainage of Makio Dam, the decision of boundary and the characteristics of divided flow are considered as the reasonable results.

Keywords: stochastic modeling, flow series, dry season, AIC

1. 研究の概要

渴水時流況を、量的・時間的に離散化して経験分布を求めるに、離散化単位の取り方にもよるが、ほぼ全体の数%にあたる個数の量の大きなデータが存在し、分布形が不連続になったり、全体の積率を左右したりしてモデル化に大きな影響を及ぼすことがよくある。そこで、本研究は流入量の大きい範囲と小さい範囲は異なる統計的性質を持つと仮定し、量の小さい範囲は自己相関性を考慮した相関同時二項分布で、量の多い範囲は無相関の分布でそれぞれモデル化を行う。ついで、両モデルの混合によって全体のモデルを表現する。全体的なモデルの適否はA I Cによって評価し最適なモデルを採択する。すなわち、渴水時流況の離散モデル化とデータの分割を同時に行う手法を示す。また、実データへの適用計算を行い、渴水時流況に対する好ましい統計モデルの選定を検討する。

2. 対象とするデータ¹⁾

量的・時間的に離散化された流入量系列 q_t ($t = 1, 2, \dots, N$) を対象とする。ここでいう量的・時間的な離散化とは、 q_t が $0, 1, \dots$ という無次元の整数値を持つことを意味する。また、 q_{t+1} は q_t に引き続き生じた流入量を意味する。ここで、

$n(i) \equiv (q_t = i)$ の発生個数, ($t = 1, 2, \dots, N$) ($i = 0, 1, \dots, c$)

$n_{ij} \equiv (q_t = i \mid q_{t+1} = j)$ の発生個数, ($t = 1, 2, \dots, N - 1$) ($i, j = 0, 1, \dots, c$)

c : 発生流入量の上限

と定義する。なお、上式に関しては、次式が成り立つ。

$$\sum_{i=0}^c n(i) = N \quad \dots \quad (1)$$

$$\sum_{i,j=0}^c n_{ij} = N - 1 \quad \dots \quad (2)$$

ここで $n(i)$ は経験分布の周辺分布に、 n_{ij} は同時分布に相当した表現である。

3. モデル化に用いる理論分布

モデル化には、周辺分布として、二項分布とポアソン分布を、相関同時分布として二項分布を用いる。これらの理論分布を以下に示す。まず、二項分布周辺分布は

$$Pb(i) = \binom{r}{i} (1-a)^{r-i} a^i \quad (i=0, 1, \dots, r) \quad \dots \dots (3)$$

$$Pb(i)=0 \quad (i=r+1, \dots)$$

また、条件付き分布は、次式となる。

$$Pb(j|i) = \sum_{s=0}^{\min(i,j)} \binom{i}{s} \binom{-i+r}{j-s} \{a(1-\rho)+\rho\}^s \\ \times \{1-a(1-\rho)\}^{s+r-i-j} a^{j-s} (1-a)^{i-s} (1-\rho)^{j+i-2s} \quad \dots \dots (4)$$

したがって同時分布 $Pbij$ は、

$$Pbij = Pb(i) \times Pb(j|i) \quad \dots \dots (5)$$

と表される。平均、分散、相関係数はそれぞれ、

$$E(q) = ra, V(q) = ra(1-a), \text{Corr}(q_t, q_{t+1}) = \rho \quad \dots \dots (6)$$

であり、パラメータ r は正整数で生起変量の上限を意味する。

つぎに、ポアソン分布周辺分布を次式で示す。

$$Pp(i) = \lambda^i \cdot \exp(-\lambda) / i! \quad \dots \dots (7)$$

4. 渇水時流況の離散化モデル¹⁾

(1) モデルの概念

流量時系列を自己相関性の違いにより、低水、高水流量集団にそれぞれ分割する手法²⁾が提案されている。本研究では、同時分布のモデルを取り扱うことで、自己相関性の有無により渴水時流況を二つの集団に分割し、自己相関性を考慮したモデル化および母集団の分割を同時に実行する。

分割の境界を $i, j = k$ ($k = 0, 1, \dots, c$) として量の少ない範囲は、一次の自己相関性を考慮した相間同時二項分布に、またそれ以外の範囲は相間のない二項、ポアソンのどちらかの分布に従うと仮定する。相間同時二項分布を用いたのは、従来、著者らが利水用貯水池問題を取り扱う際に、渴水期における貯水池系への入力として、この分布を用いた解析が有用性をもつ³⁾ことによる。まず、二つの領域で別個に求めたモデルを合成して全体のモデルを構成する（図-1 参照）。モデルを合成する際、量の多い範囲と少ない範囲とは異質の集団であり、これらが時間的に連続して生起することはないというモデルである。したがってこのモデルでは、量の多い範囲と少ない範囲との同時確率はほとんど0に等しい。もう一つは量の多い範囲と少ない範団とが時間的に連続して生起するというモデルで、両者の同時確率は合成された周辺分布の積で与える。

合成されたモデルを $MODELm(k)$ と表記する。ここで、 k は分割の境界、 m はモデル番号で表-1に示すような分布の組合せを意味する。たとえば、 $MODEL1(0)$ 、 $MODEL3(0)$ は全てのデータが無相間のポアソン分布に従うとするモデルであり、モデル番号に関係せず $MODELm(c)$ は全データが相間二項分布に従うモデルである。

(2) パラメータの推定

先述の仮定により、量の少ない範囲のモデル、つまり同時、周辺二項分布のパラメータ推定には、データとして、

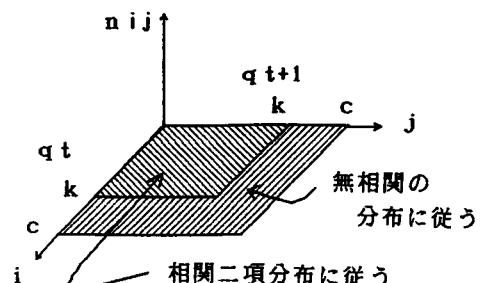


図-1 モデルの概念図

表-1 用いたモデルと分布の組合せ

	量の少ない範囲	量の多い範団	両範囲の連続生起
MODEL1	相間二項分布	周辺二項分布	しない
MODEL2	相間二項分布	ポアソン分布	しない
MODEL3	相間二項分布	周辺二項分布	する
MODEL4	相間二項分布	ポアソン分布	する

経験分布 $n(i)$ ($i = 0, 1, \dots, k$) および n_{ij} ($i, j = 0, 1, \dots, k$) を用いる。また、量の多い範囲、つまり周辺二項、ポアソン分布のパラメータ推定にはデータとして経験分布 $n(i)$ ($i = k+1, \dots, c$) を用いる。なお各々範囲外（たとえば、量の多い範囲では、 $i = 0, 1, \dots, k$ ）のデータは $n(i) = 0$ とする。

パラメータの推定方法は、積率法、最尤法などがあり一概にどの方法がよいかは判断しかねるが、著者らが数値実験的に検討した結果⁴⁾二項分布のパラメータ推定において最尤法が積率法に比べて有効であったので、ここでは最尤法を用いる。なお、二項分布は量の多い範囲と少ない範囲の双方でパラメータ推定を行うので、それらを区別するため、量の少ない範囲のパラメータに添字1 (r_1, a_1, p_1) を量の多い範囲のパラメータに添字2 (r_2, a_2) を付す。

(3) モデルの混合

ここでは、相関同時二項分布と独立な同時分布としての二項またはポアソンの両分布を合成する。まず量の多い範囲と少ない範囲が時間的に連続して生起しないとするモデルでは、たとえば、モデル番号2の場合、両範囲の同時分布を合成して、次式で表される。

$$p_{qij} = \alpha \times P_{bij} + (1 - \alpha) \times Pp(i) \times Pp(j) \quad (i, j = 0, 1, \dots, c) \quad \dots \dots (8)$$

上式において、 p_{qij} は合成された同時分布を、また α は全データの内相関二項分布に従うデータ数の割合を意味しており、本研究では次式で近似しておく。

$$\alpha = NL / (N - 1), \quad NL = \sum_{i,j=0}^k n_{ij} \quad \dots \dots (9)$$

いま、ポアソン分布のパラメータ推定にはデータとして、経験分布 $n(i)$ ($i = k+1, \dots, c$) を用いるので、推定された分布において $Pp(j) \neq 0$ ($j = 0, 1, \dots, k$) となり式(8)において、量の多い範囲と少ない範囲との同時確率 p_{qij} (i または $j = k+1, \dots, c$) はほとんど0に等しくなる。

量の多い範囲と少ない範囲が、時間的に連続して生起するとするモデルは、たとえばモデル番号4の場合、まず周辺分布の合成を行い、ついで相関二項分布と周辺分布の積を合成して次式のように表される。

$$p_{q}(i) = \alpha \times P_{b}(i) + (1 - \alpha) \times Pp(i) \quad (i = 0, 1, \dots, c)$$

$$p_{qij} = \alpha \times P_{bij} + (1 - \alpha) \times p_{q}(i) \times p_{q}(j) \quad (i, j = 0, 1, \dots, c) \quad \dots \dots (10)$$

上式で、 $p_{q}(i)$ は合成された周辺分布を意味する。また α は式(9)の値を用いる。

(4) AICの計算⁵⁾

MODEL m (k)のAICを $AIC_m(k)$ と表記し、これを計算する。 $AIC_m(k)$ は次式で定義される。

$$AIC_m(k) = -2 \times LL_m(k) + 2 \times \text{自由パラメータ数}, \quad \dots \dots (11)$$

$$LL_m(k) = \sum_{i,j=0}^c n_{ij} \times \ln(p_{qij})$$

自由パラメータ数は、相関二項分布で3個、周辺二項分布で2個、ポアソン分布で1個、ついで合成比率が1個であるから、一般的に MODEL1(k)、MODEL3(k) は6個、MODEL2(k)、MODEL4(k) は5個である。特別な場合として、全データが相関二項分布に従うとする MODEL m (c) では3個、全データが周辺二項分布に従うとする MODEL1(0)、MODEL3(0) では2個、全データがポアソン分布に従うとする MODEL2(0)、MODEL4(0) では1個となる。

各モデル番号、各分割のうちで AIC が最小のモデルを最適なモデルとして採択する。一般に複雑な混合をすればするほど、自由パラメータ数の増加、すなわち AIC の増大に連なり、そのようなモデルは採択されにくくなる。しかしそのようなモデルがあえて採択された場合には、流量時系列の母集団が統計的特性の異なる複数の集団から形成されていたといえるだろう。普通あまり自由パラメータ数の多いモデルはパラメータ推定が不安定になるなどの理由で分布としては適切でないが、本研究では、AIC だけで、モデルの適否を判断するため、自由パラメータ数が5個や6個といったモデルを採択することも起りえる。

5. 模擬発生データによる検討

表2 最適モデルの選定例 ($r_1 = 5, a_1 = 0.4, \rho = 0.6, r_2 = 10, a_2 = 0.8, \alpha = 0.9, N = 500$)

パラメータの再現性を検討するため、模擬データに対し、最適なモデルを策定する。量の少ない範囲の母集団として、 $r = 5, a = 0.4, \rho = 0.6$ の相関二項分布を、量の多い範囲の母集団として $r_2 = 10, a_2 = 0.8, \alpha = 0.9$ の周辺二項分布を仮定した。

それぞれの範囲の分布を、混合の比率 $\alpha = 0.90, 0.95, 0.99$ の3種で合成し、全体のモデルを構成した。合成は、量の多い範囲と少ない範囲とは時間的に連続して生起しない、つまり先述のMODEL1に従う形で行った。合成された分布に従う、全データ数 $N = 100, 500, 1000$ 個の頻度分布を模擬データとし、MODEL1およびMODEL2によってモデル化した。二項分布の確率は、式(6)で表現されるから、量の少ない範囲の母集団の平均は、2.0、量の多い範囲の母集団の平均は、それぞれ5.0, 6.0, 7.0, 8.0, 9.0となり、パラメータ a_2 が小さいほど両母集団の平均値に差がなくなり、データの分割が難しくなることが予想される。

最適モデルの策定例として、量の多い範囲の母集団として、 $r = 10, a = 0.8$ の周辺二項分布を合成の比率 $\alpha = 0.9$ で合成した分布に従うデータ数500個の離散分布に対して、MODEL1, MODEL2でそれぞれモデル化した結果を表-2に示す。表の中でNoは、分割の境界、すなわち3. (1)で述べた k に対応する。まず、全体で AIC を比較すると、MODEL1のNo5がAICが最小なモデルとなっている。2番目にAICが最小なモデルであるMODEL1のNo6とのAICの差は、約9あるので、このモデルが全体の中で、最適なモデルといって良いであろう。合成される前の母集団のパラメータや、合成の比率もほぼ正確に推定されており、パラメータの再現性も非常に良好である。MODEL2の中でAICが最小なのは、No5であり分割の境界は正しく推定されている。これは、主体となる分布（この場合相関同時二項分布）のモデルが適切であれば、量の多い範囲のモデルがそれほど適切でなくとも、分割の境界は正しく推定できることを意味していると思われる。

つぎに、データ個数 N や混合率 α 、さらに量の多い範囲の母集団の統計的性質がパラメータ推定に与える影響を考察する。 $N = 100, 500, 1000$ の場合で、 α やパラメータ a_2 とパラメータ推定精度との関係を表3～5に示す。表中で \times は両方の分布でのパラメータ推定の失敗を、 \triangle は量の少ない範囲のパラメータ推定の成功を、 \circ は両方の分布でのパラメータ推定の成功を意味する。パラメータ推定の成否は、推定誤差が $\pm 5\%$ 以下であれば、成功したと判断した。まず $N = 100$ の場合は、 α や a_2 に関係なくほとんどパラメータ推定に失敗している。本研究では、比較的の自由パラメータ数の多いモデルを考えているため、100個というデータ個数は少ないとと思われる。 $N = 500$ の場合は、 α が一定の場合 a_2 が大きくなるほどパラメータの推定精度はよ

No	α	r_1	a_1	Corr.	MODEL 1			MODEL 2	
					r_2	a_2	AIC1	λ	AIC2
0	0.	—	—	—	10	0.2565	4422.5034	2.5648	3995.3823
1	0.3077	1	0.7040	0.3140	10	0.3363	4309.3496	3.3628	3962.0828
2	0.6215	2	0.6590	0.3980	10	0.4492	4034.4795	4.4920	3750.4202
3	0.8300	3	0.5836	0.4828	24	0.2634	3405.2039	6.3214	3383.2705
4	0.8988	4	0.4880	0.5580	10	0.7900	3118.4221	7.9000	3164.8184
5	0.9069	5	0.3981	0.6100	10	0.8152	3051.6685	8.1522	3148.0219
6	0.9130	6	0.3321	0.6390	10	0.8302	3060.2493	8.3023	3156.3796
7	0.9332	7	0.2904	0.6692	10	0.8697	3092.9851	8.6970	3173.5242
8	0.9636	25	0.0876	0.7520	10	0.9278	3217.4058	9.2778	3247.6682
9	0.9899	25	0.0978	0.7710	10	1.0000	3412.1382	10.0000	3335.8369
10	1.0000	25	0.1026	0.7760	—	—	3411.8601	—	3411.8601

表-3 パラメータ推定精度の比較($N=100$)

α	a_2	0.5	0.6	0.7	0.8	0.9
0.90	\times	\times	\times	\times	\triangle	
0.95	\triangle	\triangle	\times	\times	\triangle	
0.99	\times	\times	\times	\times	\times	

表-4 パラメータ推定精度の比較($N=500$)

α	a_2	0.5	0.6	0.7	0.8	0.9
0.90	\times	\triangle	\triangle	\circ	\circ	
0.95	\times	\triangle	\triangle	\circ	\circ	
0.99	\triangle	\triangle	\triangle	\triangle	\circ	

表-5 パラメータ推定精度の比較($N=1,000$)

α	a_2	0.5	0.6	0.7	0.8	0.9
0.90	\times	\times	\circ	\circ	\circ	
0.95	\times	\triangle	\triangle	\circ	\circ	
0.99	\times	\triangle	\triangle	\triangle	\circ	

くなっている。これは、混合率が等しいときは、両分布の統計的性質が異なるほど推定精度（特に量の多い範囲のパラメータ推定）がよくなることを意味する。また、 a_2 が一定の場合、 a_2 が小さい、すなわち両範囲の平均が近い場合は、 α が大きくなるほど、推定精度がよくなる。逆に a_2 が大きい、すなわち両範囲の平均が離れている場合には、 α が大きくなるほど推定精度は悪くなる。両範囲の平均が近いと、 α が大きくなることにより量の多い範囲のデータ個数が減少し、量の少ない範団のパラメータ推定に雑音の減少として推定精度がよくなるような影響を与えるためと思われる。

しかし、データの分離が行われ難いことにより、量の多い範団のパラメータ推定は、ほとんど期待できない。両者の平均が離れていると、データの分割は行われやすく、量の少ない範団のパラメータ推定はデータ個数が多いこともあり、ほぼ正確に行われる。また α が大きくなると、量の多い範団のデータ個数の減少に連なり、パラメータの推定精度が悪くなる。 $N = 1,000$ の場合、全体的な傾向は $N = 500$ の場合とほとんど同じで、その傾向がより顕著に現れているようにみえる。

これら、模擬データによる検討により、本研究で示した方法でデータの分離、モデル化をする場合には、少なくとも数百個のデータと、ある程度、両母集団の統計的性質が異なっていることが必要だと思われる。しかし、たとえパラメータが正確に再現されなかったとしても、得られたモデルは、データ個数との兼ね合からみてデータをよく表現しているモデルと考えるべきであろう。

6. 実データに対する適用計算

牧尾ダムの1969年～1986年冬期渇水期（12月1日から翌年2月28日まで）日流入量に対し、適用計算を行った。データの組は、各水期ごと（全17期間分）および全渇水期を統合したもの、の18組であり、各組を、単

表-6 最適モデルの選定例（1969～1986年、牧尾ダムの冬期渇水期流入量）

位期間1, 3, 5日、 単位量=単 位期間×3 m^3/sec に より離散化 し、MODEL1 ～MODEL4で モデル化を 行い、最適 なモデルを 選定した。	No	α	r1	a1	Corr	MODEL 1, 3				MODEL 2, 4		
						r2	a2	AIC1	AIC3	λ	AIC2	AIC4
0	0.	-	-	-	13	0.1164	2076.768	2076.768	1.5131	1975.198	1975.198	
1	0.6961	1	0.7958	0.7352	13	0.2399	1741.075	1741.043	3.1183	1658.080	1680.177	
2	0.8824	2	0.5090	0.7830	13	0.3761	1756.075	1593.677	4.8889	1659.270	1552.288	
3	0.9346	3	0.3810	0.7550	13	0.4923	1778.161	1581.625	6.4000	1673.342	1557.093	
4	0.9510	5	0.2418	0.7620	50	0.1440	1698.496	1587.996	7.2000	1684.377	1587.183	
5	0.9673	10	0.1270	0.7460	31	0.2677	1737.981	1637.154	8.3000	1715.178	1637.167	
6	0.9771	50	0.0263	0.7322	21	0.4422	1769.815	1685.871	9.2857	1738.184	1686.718	
7	0.9837	50	0.0270	0.7038	17	0.6000	1799.182	1716.085	10.2000	1756.012	1718.235	
8	0.9902	9	0.1547	0.6342	13	0.8974	1941.289	1797.502	11.6667	1801.806	1804.432	
9	0.9902	10	0.1391	0.6400	13	0.8974	1920.020	1788.121	11.6667	1797.990	1795.051	
10	0.9935	11	0.1289	0.5880	13	0.9615	1965.420	1853.074	12.5000	1816.760	1862.964	
11	0.993	12	0.1182	0.5912	13	0.9615	1958.409	1845.213	12.5000	1814.104	1855.104	
12	0.9967	13	0.1129	0.5751	13	1.0000	1983.689	1894.810	13.0000	1829.807	1902.636	
13	1.0000	14	0.1082	0.5380	-	-	1963.041	1963.041	-	1963.041	1963.041	

選定例と

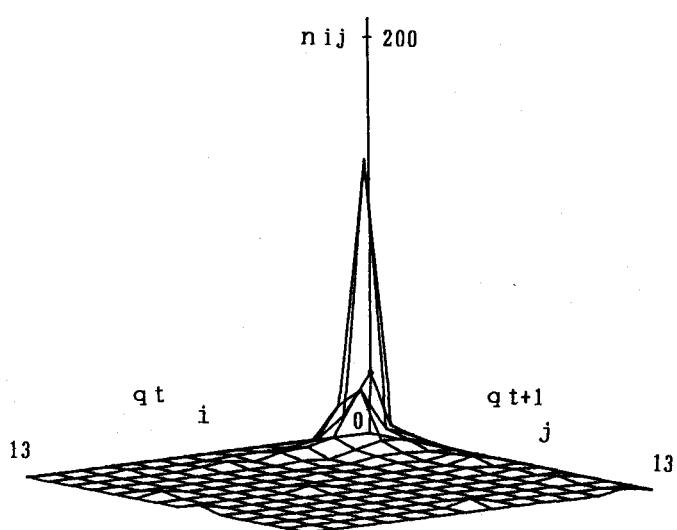


図-2 牧尾ダム流域での渇水期流量の同時頻度分布

して、全渴水期の日流量を、単位期間5日で離散化したデータの同時頻度分布を図-2に、最適モデルの選定結果を表-6に示す。図-2に示した頻度分布をみると、量の大きなデータが量の少ないデータの集団から、飛び離れて存在していることが解る。

さてこれら量の多いデータを量の少ないデータから分離したモデル化を行こなう。表-6のモデル選定結果をみると、全体では、MODEL4(2)がAICが最小で4つのモデルのうち、最適なモデルであることが解る。混合率 α は0.88、自己相関係数 $\rho = 0.78$ であるから、全体の9割が自己相関性の強い流量時系列成分から構成されていることが推測される。2番目にAICが小さいのは、やはり、MODEL4(3)であるから、モデルとしては、MODEL4が適しているようだ。分布形は同じ組み合せで、量の多い範囲と少ない範囲とが時間的に連続して生起しないとするMODEL2のうちでAICが最小なのは、MODEL2(1)だが、MODEL4(2)とのAICの差は、100以上もある。また、MODEL1とMODEL3のAICを比較すると、全体的にみて、MODEL3のAICの方が小さくなっている。これらの結果から、このデータの場合は、量の多い範囲と少ない範囲が時間的に連続して生起するとしたモデルの方が妥当なように思える。

ついで、各渴水期ごとに最適なモデルを求め、単位期間ごとに最適モデルの内訳を表-7に示す。表中でMODEL(c)というのはモデル番号に関係せず、全データが相関二項分布に従うとするモデルである。まず、全体的にみて、MODEL(c)は数ケースしかないので、今回のデータの場合は、量の多い範囲と少ない範囲とにデータを分割するモデルの方が適していたようだ。また、単位期間1日の場合は、かなり平均的にどのモデルも採用されているが、単位期間3日、5日ではMODEL3が最も多く採用されているので、全体的にみても、量の少ない範囲と多い範囲が連続して生起するというモデルの方が妥当なようである。

なお、このようなモデルの使用方法としては、たとえば、分割されたデータのうち、自己相関の強い流量時系列成分だけを貯水池系への渴水時入力として採用するのが実用的であろう。

7. まとめ

本研究では、渴水時流況データを流量時系列成分とそれ以外との成分に分離してモデル化する手法を示した。モデルとして、相関同時二項分布と、周辺二項分布、ポアソン分布の二つの分布を合成したものを用いた。しかし、基本的に離散分布であれば、どんな分布からモデルを構成してもよい。また、モデルを構成する分布の数も二つにかぎられるわけではない。それぞれのデータに適した分布を用いることで、データの分割や、統計的特性の単純化をより正確に行うことができると思われる。しかし、採用分布が異なった場合でも、本研究で示した手法で、AIC基準によった最適モデルの選定が可能となろう。

8. 参考文献

- 1) 鈴木正人・長尾正志：AICを用いた相関渴水時流況の離散モデル化の研究，第44回年次学術講演会概要集II部門，1989
- 2) 室田明・神田徹：利水を対象とした流量時系列の解析について，水理講演会講演集，1989
- 3) 鈴木正人・長尾正志：2段階推移モデルによる相関離散分布を受ける貯水池理論，土木学会論文集，第411号／II-12, pp.161-168, 1989
- 4) 鈴木正人・外山康彦・長尾正志：渴水時流況の表現のための離散化確率モデルの選定，土木学会中部支部講演概要集，1988
- 5) 坂元慶行・石黒真木夫・北川源四郎：情報量統計学，共立出版社，1983