

データマイニングにおけるGA・情報エントロピーの応用

Data Mining using Information Entropy and Genetic Algorithm

須藤 敦史*・渋谷 卓**・星谷 勝***

Atsushi SUTOH, Taku SHIBUYA and Masaru HOSHIYA

* 博士(工学) (株) 地崎工業 土木技術部 主席研究員(〒105-8488 東京都港区西新橋2-23-1)

** (株) 地崎工業 情報システム部 (〒105-8488 東京都港区西新橋2-23-1)

*** Ph.D. 武藏工業大学教授 工学部土木工学科 (〒158-8557 東京都世田谷区玉堤1-28-1)

In this study consists of the following two topics, one is a basic consideration on a data mining which is the power of current data-processing functions, of interesting knowledge rules from huge database. And the other, we introduce data mining with a view to discuss applications of artificial life and decision tree procedures. Data mining procedures which, decision tree using information entropy theory and genetic algorithm are proposed, and red tide data from Tokyo bay were analyzed. Finally, it is found that the usefulness of these data mining procedures for non-structural system identification.

Key Words: data mining, decision tree, genetic algorithm, information entropy, identification

1. はじめに

現在,データウェアハウス活用の要求が多くなっているが定性・定量的データが混在し,加えて相互関係が複雑であるため,現状の解析技術では有効利用がなされていない.そこで価値ある情報や知識の発掘を目的としたデータベースからの知識発見 (KDD:Knowledge Discovery in Databases) あるいはデータマイニング (DM:Data Mining)^{①,②}が注目されている.

もっとも,膨大なデータから知識を得ようとする研究は確率・統計,機械学習など多様な枠組みで試みられている^③が,データマイニングは図-1 に示すようにデータ中の隠れたルールを客観的に発見するための現実問題として従来の手法を統合した新たなプロセッシング技術あるいは考え方であり,特にマーケティング分野で多くの応用例が報告されている^{④,⑤}.

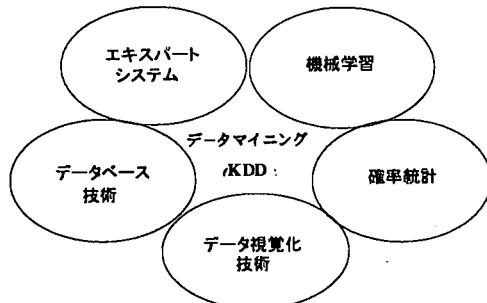


図-1 データマイニングの概念図

本研究は,事象の相関関係(決定木)に条件付き情報エントロピーを用いた評価方法を提案し,同時にGAによ

り東京湾で観測された水質調査データと赤潮発生との相関解析を行っている.

2. データマイニング

データマイニングの背景には以下の要因が影響していると考えられる.

(a)データベースの発展

コンピュータの進歩で膨大かつ多様なデータの蓄積が進んでおり,これらの有効活用が求められている.

(b)理論・技術の統合

従来のデータ解析理論もしくは技術を統合・システム化した新しい解析手法が求められている.

(c)技術のソフトウェア化

実データを解析するため汎用化され,かつ操作手順が簡単なツールが求められている.

一般的にデータマイニングは「仮説検証型」と「仮説生成型」に分けられ,前者は仮説のデータによる検証すを,後者はデータを単純な表現形式に変換して隠れたルールの発見することを目的としている.しかし,両者とも従来のデータ解析に比べて「状況予測」より「結果解釈」に重点を置いているところが特徴である.

3. データマイニングに用いた解析ツール

基本的にはデータ解析技術・ツールならば何を用いても構わないが目的に応じて様々なツールを単独あるいは複数組合せて使用する.本研究では以下の2手法により

データマイニングを行っている。

(a) デシジョンツリー

母集団を属性ごとに分割し、木の枝のように表現する手法であり「決定木」あるいは「判別木」とも呼ばれる。この手法は複数のルールを同時に表現できるため、事象全体を把握するのに有効な手法である。

一方、情報エントロピーは事象の不確定の度合いを表しているため、事象間の相関強さを定量的に評価できる利点がある。加えて相互情報量は条件付き確率と解釈できるため、事象の時間的な前後関係を明確にする可能も有している。

そこで、本研究では決定木における各事象（属性）の相関関係の評価指標としての情報エントロピーや相互情報量および相関関係の整理を行っている。

(b) 遺伝的アルゴリズム¹⁾

生物の進化過程・集団遺伝のメカニズムを工学モデル適用にした手法であり、主に組み合せ最適化問題などに応用されている。本研究では事象のランダムな組み合せを効率よく探索するために利用している。

4. 東京湾における赤潮と観測項目との相関関係

(1) 水質観測項目

観測データは図-2に示す東京湾（西側）の10点で観測された水質調査データと赤潮発生の有無を用いており、項目の詳細を表-1に示す。

ここで観測データは1988~92年の5年間における4~9月に観測されたものであり、データ総数は300である。

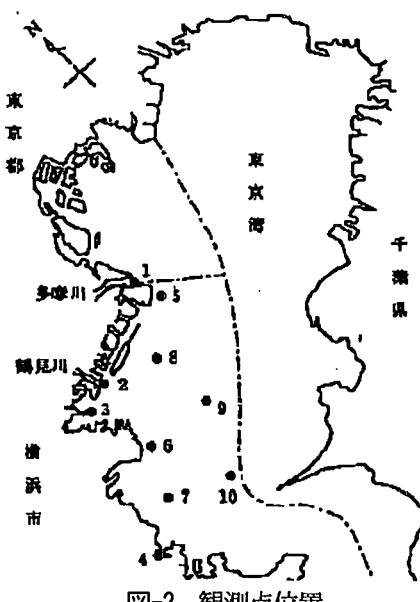


図-2 観測点位置

として赤潮発生と観測項目との相関解析を行う。

表-1 観測項目

	説明	
A	気温	
B	水温	
C	透明	
D	pH	
E	COD	化学的酸素要求量(mg/l)
F	DO	溶存酸素量(mg/l)
G	T-P	全リン(mg/l)
H	PO4-P	リソ酸態リソ(mg/l)
I	T-N	全窒素(mg/l)
J	NH4-N	アノニア態窒素(mg/l)
K	NO2-N	亜硝酸態窒素(mg/l)
L	NO3-N	硝酸態窒素(mg/l)
M	SAL	塩分(mg/l)
N	Chl-a	クロフィルa(mg/l)
O	赤潮	赤潮発生の有無

表-2 データのプール属性化

	平均値未満	平均値以上
A	A ₁	A ₂
B	B ₁	B ₂
C	C ₁	C ₂
D	D ₁	D ₂
E	E ₁	E ₂
F	F ₁	F ₂
G	G ₁	G ₂
H	H ₁	H ₂
I	I ₁	I ₂
J	J ₁	J ₂
K	K ₁	K ₂
L	L ₁	L ₂
M	M ₁	M ₂
N	N ₁	N ₂

※ 赤潮あり:O₁ 赤潮なし:O₂

(2) 決定木（デシジョンツリー）による分析

(a) 情報エントロピー

情報エントロピー¹⁾は「現象や情報の不確定の度合い」を表す尺度であり、事象間の相関強さを評価することが可能となる。

いま、離散型確率分布を有する事象Xを考える。

$$X = \begin{pmatrix} X_1 & X_2 & \cdots & X_m \\ P_1 & P_2 & \cdots & P_m \end{pmatrix} \quad (1)$$

ただし、 $0 \leq p_i \leq 1$ かつ $\sum p_i = 1$
このとき事象Xの情報エントロピーは式(2)となる。

$$H(X) = H(p_1, \dots, p_m) = - \sum_{i=1}^m p_i \log p_i \quad (2)$$

ここで観測項目のA~Nは数値属性、赤潮に関する項目Oはプール属性(0-1)のデータであるが、簡略化するため表-2のように全データを平均値以上・以下のプール属性化し、それらを条件とした「If ~then ...」形式のルール

ここで式(3)に示す確率分布を有する事象 Y を考える。

$$Y = \begin{pmatrix} Y_1 & Y_2 & \cdots & Y_n \\ q_1 & q_2 & \cdots & q_n \end{pmatrix} \quad (3)$$

ただし、 $0 \leq q_j \leq 1$ かつ $\sum q_j = 1$

また Y_j が条件として与えられたときの X_i の条件付きエントロピーは式(4)で与えられ、範囲は式(5)となる。

$$H(X|Y) = \sum_{j=1}^n q_j H(X|Y_j) \quad (4)$$

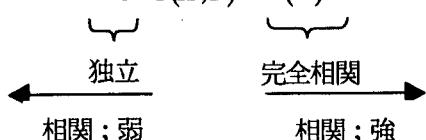
$$0 \leq H(X|Y) \leq H(X) \quad (5)$$

ここで条件が与えられたことによる条件付き情報エントロピーの差は相互情報量と呼ばれ、式(6)となる。

$$I(X;Y) = H(X) - H(X|Y) \quad (6)$$

相互情報量 $I(X;Y)$ は事象 X と Y の相関の強さを表す指標として用いることができ、加えて事象の時間的な前後関係を明確にすることが可能となる。

$$0 \leq I(X;Y) \leq H(X) \quad (7)$$



(b) 事象（観測項目）の相互情報量

各事象に対する条件付き確率の差、相互情報量に相関係数を加えた比較表を表-3に示す。

表-3 赤潮との相間関係

X	$ P(O_1 X_1) - P(O_1 X_2) $	I(O;X)	R(O;X)
A 気温	0.002	0.000	0.023
B 水温	0.041	0.002	0.021
C 透明	0.265	0.085	-0.382
D pH	0.176	0.032	0.303
E COD	0.321	0.106	0.603
F DO	0.194	0.036	0.467
G T-P	0.163	0.023	0.222
H PO4-P	0.096	0.008	-0.141
I T-N	0.080	0.005	0.036
J NH4-N	0.143	0.020	-0.161
K NO2-N	0.058	0.003	-0.076
L NO3-N	0.125	0.014	-0.139
M SAL	0.025	0.001	0.050
N Chl-a	0.563	0.199	0.588

表より、条件付き確率の差、相互情報量は事象の相関の強さを示しており、比較表から絶対値で比較するとほぼ同じ傾向を示している。

ここで表-3より相関の強い順に分割した決定木1を図-3に示す。決定木における分割終了条件は「確信度が0 or 1」もしくは「サポートが0.15以下」まで、図中の各属性（ノード）の数値を以下に示す。

一段目： $m_1 (m_1, m_2)$ 条件を満足する総数： m_1 、その中で赤潮発生数： m_1 、未発生数： m_2

二段目：確信度 m_1/m (赤潮の発生数/条件を満足する数) で定義される指標で条件付き赤潮発生確率

三段目：サポート $m_1/300$ (条件を満足する数/データ総数) で定義される指標で現段階まで条件を満足する確率

四段目：全赤潮発生数に対する割合 $m_1/47$ (現段階の赤潮発生数 m_1 /総数中の全赤潮発生数 (47)) 詳細な相関関係評価のための指標として新たに定義

(c) 決定木による解析結果

従来では、確信度が高いほどルールとしての評価は高いが、決定木1（図-3）において矢印が太い関係に着目するとノード分岐が進むに従い確信度は増加傾向にあるが、全赤潮発生数に対する割合は逆に減少している。特に「DO：大」という条件が追加するとその減少が大きいため、条件を満足するデータを分母とする確信度の客観性に疑問が残る。

そこで、確信度と新たに定義した全赤潮発生数に対する割合が大きい事象までを示すと「Chl-a：大 \cap COD：大 \cap 透明：小」ならば「赤潮発生」となり、これは全赤潮発生数に対する割合は0.872 (41/47)と高く、かつ確信度も0.651 (41/63)と高い。よって両指標が高い値を示す事象までを「1次的原因」として考える。

次に確信度のみが高い事象は「Chl-a：大 \cap COD：大 \cap 透明：小 \cap DO：大 \cap pH：大 \cap T-P：大」ならば「赤潮発生」となる。この確信度は0.759 (22/29)と高いが、全赤潮発生数に対する割合は0.468 (22/47)と0.5を下回った値を示している。つまり赤潮が発生している47ケースの中で半分以上のデータがこの条件を満足していない。そこで確信度は増加しているが全赤潮発生数に対する割合が減少してしまう「2次的原因」として考える。

よって決定木1から導かれる相関の強さ（ルール）を上記の評価値から分類すると以下のようになる。

1次的原因：「Chl-a：大 \cap COD：大 \cap 透明：小」

2次的原因：「DO：大 \cap pH：大 \cap T-P：大」

次に、条件付き確率の差および相互情報量を参考にして分割した決定木2を図-4に示す。

決定木2ではツリーの構造が簡素化されているため、事象の相互関係が複雑なものに対して、効率的な分割を行えると考えられる。

ここで図-4から得られる赤潮発生との相関の強さを同様に評価値から分類すると以下のようになる。

1次的原因：「Chl-a：大 \cap 透明：小 \cap COD：大」

2次的原因：「水温：小」

1次的原因は決定木1と同様の結果、また2次的原因は「水温：小」という事象が追加されてる。

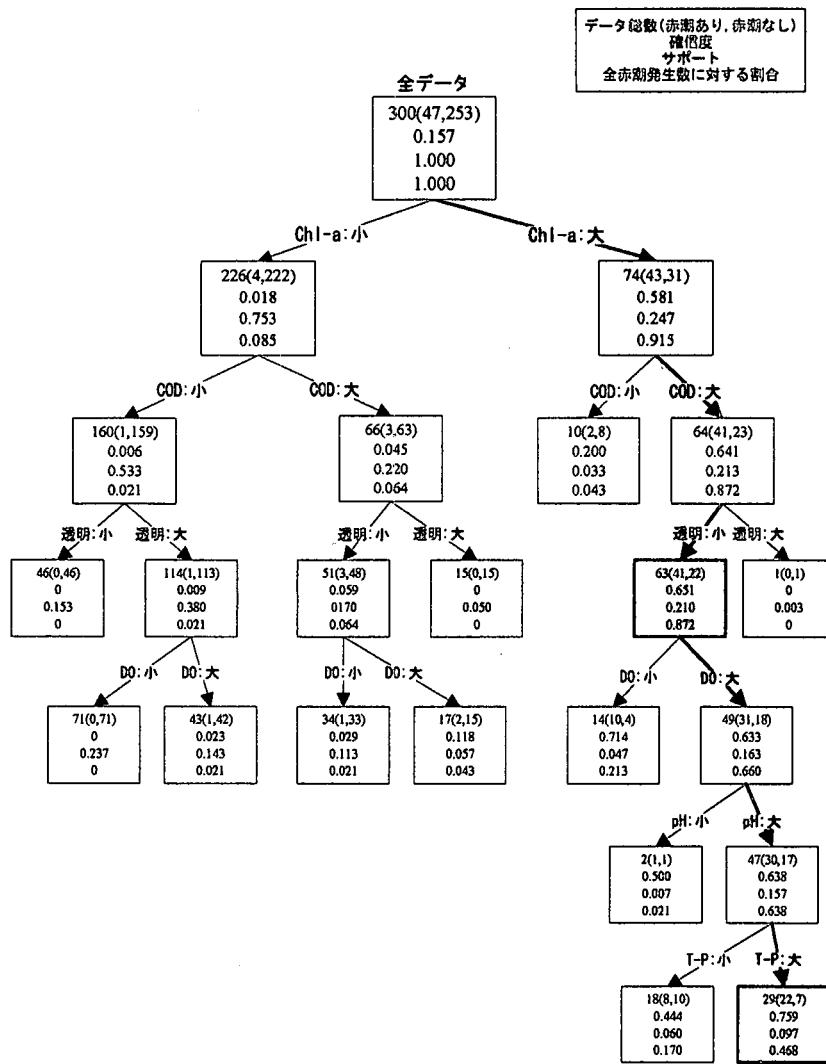


図-3 決定木 1

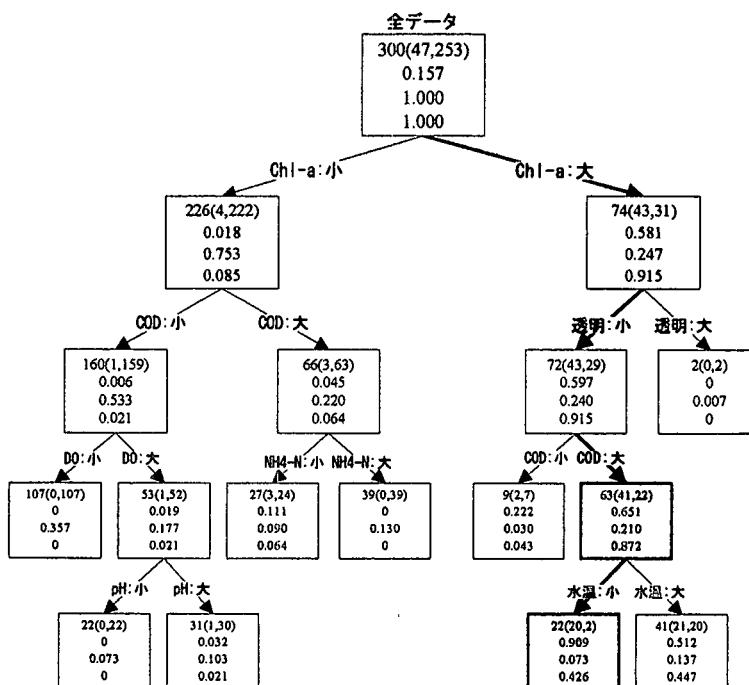


図-4 決定木 2

したがって、決定木1では獲得できなかった「水温：小」が条件付き確率の差もしくは相互情報量を参考にするにより抽出できた。

(2) GAによる分析

決定木の特徴は、条件付き確率の差もしくは相互情報量を指標として指標の大きさ（相関の強い）順に枝を形成していく点にある。しかし一般的にデータベースは数多くの事象（項目）を有し、データ総数も膨大であるため、各事象における条件付き確率の差や相互情報量を求めるには多くの計算時間が必要となる。

そこで相関の順序に関係なく事象の組合せを検討するためにGAを用いたによる相関の強さ（ルール）解析を行う。

(a) GAの解析手順

GAにおける目的関数を「①確信度」、「②全赤潮発生数に対する割合」として、これを最大にする事象の組み合わせを探査する最適化問題として、赤潮と事象との相関関係（ルール）を求める。

ここで組み合わせる属性の数は3, 4, 5と限定し、制約条件としてサポートが0.15以下の属性の組み合せは削除している。またGAにおける世代数300、個体数500、交叉確率0.9、突然変異確率0.5を設定した。

(b) GAの解析結果

組み合わせる事象数が3つの結果を表-4に示す。表中の「1」は選択された事象であり、A_i～N_j, m, n_iは決定木と同様である。また、事象数が4つの結果を表-5、5つの結果を表-6に示す。

以下に各目的関数の高い組み合せを挙げる。ここで（）内は確信度、全赤潮発生数に対する割合である。

- (ア) 3事象-目的関数① (0.71429, 0.74468) 「COD : 大 \cap NH4N : 小 \cap Chl-a : 大」
(イ) 4事象-目的関数① (0.72917, 0.74468) 「COD : 大 \cap NH4N : 小 \cap Chl-a : 大 \cap 透明 : 小」
(ウ) 5属性-目的関数① (0.72340, 0.72340) 「透明 : 小 \cap COD : 大 \cap Chl-a : 大 \cap NO3-N : 小 \cap pH : 大」
(エ) 3事象-目的関数② (0.65079, 0.87234) 「透明 : 小 \cap COD : 大 \cap Chl-a : 大」
(オ) 4事象-目的関数② (0.72000, 0.76596) 「COD : 大 \cap pH : 大 \cap Chl-a : 大 \cap 透明 : 小」
(カ) 5事象-目的関数② (0.72340, 0.72340) 「COD : 大 \cap Chl-a : 大 \cap 透明 : 小 \cap NO3-N : 小 \cap pH : 大」

以上より「確信度」、「全赤潮発生数に対する割合」のどちらを目的関数に用いてもほぼ同じ事象が選定される。ここで、下線はどちらかの目的関数により得られた属

性である。

- また目的関数①, ②により得られた事象を以下に示す。
①, ②共通：「Chl-a : 大 \cap COD : 大 \cap 透明 : 小」
②：「pH : 大 \cap NO3-N : 小」, ①：「NH4N : 小」

5. 結論

本研究ではデータマイニングにおける事象間の相関関係が条件付き確率の差および相互情報量によって評価できることを示し、同時に東京湾で観測された水質観測データを用いて赤潮発生要因の相関解析を条件付き確率の差および相互情報量を用いた決定木の分類および膨大な組み合わせのためにGAを用いて行い、以下の結論が得られた。

- (1) 条件付き確率では「Chl-a」, 「COD」, 「透明」の事象が抽出され、また相互情報量を用いて評価した場合も同様の結果が得られた。
- (2) 条件中の割合の評価により、2次的な要因「DO : 大 \cap pH : 大 \cap T-P : 大」や「水温 : 小」が赤潮発生に対して相関を有する結果となった。
- (3) GAによる解析からも「Chl-a : 大 \cap COD : 大 \cap 透明 : 小」が選択された。また決定木では得られなかった「pH : 大 \cap NO3-N : 小」, 「NH4N : 小」が選択された。

また、今後の課題としては以下の点が挙げられる。

- (1) データの前処理が重要となるため、効率的な数値データの前処理について検討が望まれる。
- (2) 「A \rightarrow B \rightarrow 赤潮発生」という因果関係などの時間的な分析が期待される。

参考文献

- 1) Pieter Adrians, Dolf Zantinge 著 山本英子・梅村恭司訳：データマイニング、共立出版、1998
- 2) 大規模データベースからの知識獲得、人工知能学会誌、Vol.12, No.4, pp496-549, 1997.7
- 3) 徳山豪：データマイニングに使われる最適化の数理、応用数理、VOL.6, NO.4, pp303-313, 1996.12
- 4) 中林三平：データマイニング 値ある情報を掘り当てる、NIKKEI COMPUTER, pp142-147, 1996.9.30.
- 5) 須藤敦史、高須光朗、星谷勝：ニューラルネットワークを用いたデータマイニングによる非構造システムの同定、応用力学論文集、Vol.2, pp.83-90, 1999.
- 6) 有本卓：確率・情報・エントロピー、森北出版、1992.
- 7) 北野宏明：遺伝的アルゴリズム、産業図書、1993.

表-4 事象：3の結果 (GA)

		条件の属性																m	m ₁	確信度	サポート	全赤潮発生数に対する割合							
		A ₁	A ₂	B ₁	B ₂	C ₁	C ₂	D ₁	D ₂	E ₁	E ₂	F ₁	F ₂	G ₁	G ₂	H ₁	H ₂	I ₁	I ₂	J ₁	J ₂	K ₁	K ₂	L ₁	L ₂	M ₁	M ₂	N ₁	N ₂
①.「確信度」による評価																													
②.「全赤潮発生数に対する割合」による評価																													
①の合計		0	0	0	0	2	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	
②の合計		0	0	0	0	10	0	0	1	0	9	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	
①+②の合計		0	0	0	0	12	0	0	1	0	17	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	

表-5 事象：4の結果 (GA)

		条件の属性																m	m ₁	確信度	サポート	全赤潮発生数に対する割合							
		A ₁	A ₂	B ₁	B ₂	C ₁	C ₂	D ₁	D ₂	E ₁	E ₂	F ₁	F ₂	G ₁	G ₂	H ₁	H ₂	I ₁	I ₂	J ₁	J ₂	K ₁	K ₂	L ₁	L ₂	M ₁	M ₂	N ₁	N ₂
①.「確信度」による評価																													
②.「全赤潮発生数に対する割合」による評価																													
①の合計		0	0	0	0	10	0	0	1	0	10	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0		
②の合計		0	0	0	0	10	0	0	7	0	8	0	0	0	0	0	0	0	0	1	0	0	0	4	0	0	0		
①+②の合計		0	0	0	0	20	0	0	8	0	18	0	0	0	0	0	0	0	0	0	0	0	0	0	2	0	0		

表-6 事象：5の結果 (GA)

		条件の属性																m	m ₁	確信度	サポート	全赤潮発生数に対する割合							
		A ₁	A ₂	B ₁	B ₂	C ₁	C ₂	D ₁	D ₂	E ₁	E ₂	F ₁	F ₂	G ₁	G ₂	H ₁	H ₂	I ₁	I ₂	J ₁	J ₂	K ₁	K ₂	L ₁	L ₂	M ₁	M ₂	N ₁	N ₂
①.「確信度」による評価																													
②.「全赤潮発生数に対する割合」による評価																													
①の合計		0	0	0	0	10	0	0	6	0	9	0	1	0	0	0	0	0	0	7	0	4	0	10	0	0	1	0	
②の合計		0	0	0	0	9	0	0	9	0	2	0	0	0	4	0	0	0	0	1	0	9	0	0	0	0	0	7	
①+②の合計		0	0	0	0	19	0	0	15	0	11	0	1	0	0	4	0	0	0	0	16	0	5	0	19	0	0	0	9