

# 航空宇宙技術研究所報告

TECHNICAL REPORT OF NATIONAL AEROSPACE LABORATORY

TR-854

大規模衛星画像クラスタ解析の研究

ヒストグラム・モード法

松 本 甲 太 郎 ・ 中 正 夫 ・ 山 本 浩 通

1985 年 3 月

航空宇宙技術研究所

NATIONAL AEROSPACE LABORATORY

# 大規模衛星画像クラスタ解析の研究\*

## ——ヒストグラム・モード法——

松 本 甲 太 郎\*\*      中      正 夫\*\*      山 本 浩 通\*\*

### A Study of the Clustering Method for a Huge Satellite Image

#### — Histogram Mode Method —

By

Kohtaro Matsumoto, Masao Naka, Hiromichi Yamamoto

#### ABSTRACT

The clustering technique is important as the unsupervised classifier in the realization of an automated satellite image analysis system. But current cluster analysis methods are not suitable to deal with huge, multi-dimensional satellite images because of the dependency on initial partitions and the tremendous computational time required.

In this paper, to analyze a huge satellite image, we will describe two new algorithms, i.e., Multi-Layered Hashing Scheme and Histogram Mode Method. The first one is for the construction of a multi-dimensional histogram from the huge satellite image data. And the second one is for deriving natural clusters from the histogram using the basic idea above.

The algorithms were implemented on the large scale pipeline computer (FACOM 230/75 APU), and have been successfully used. Some results are shown in the paper.

#### は    じ    め    に

クラスタ解析手法は衛星画像の自動解析の教師無し分類手法として重要であるが、従来のクラスタ解析手法では少数のデータしか取り扱うことが出来なかったため、膨大な衛星画像の極く小さな部分画像の解析が可能であるにすぎなかった。このため、部分画像の選定が問題となり、教師無し分類としてのクラスタ解析のメリットが十分発揮できなかった。本論文はこのような点から、衛星画像のフレーム全体を1度に解析できる高速なクラスタ解析手法について論ずる。

広義のクラスタ解析は“似た者同士を集める”統計学的手法であり、解析対象データから単純に似た者同士をまとめる場合と、解析対象データからある特徴をもった少数の代表的クラスタを抽出する場合がある(文献1)。リモートセンシングにおいて衛星画像の自動解析に用いるクラスタ解析は、主に後者の場合である。解析対象(衛星画像)から地上のカテゴリーに応じた特徴をもったクラスタを抽出し、これによって衛星画像を弁別分類し地上の被覆分類図を作成する事が衛星画像クラスタ解析の目的である。

衛星画像のクラスタ解析には、これまで大別して2種類の手法が使われて来た。1つはBall & Hallの非階層的再配置法の流れをくむもので代表的なブ

\* 昭和60年1月9日受付

\*\* 計算センター

ログラムに ISODATA, ISOCLS, LARSYS (文献 12, 5, 20) 等があり, 手軽で高速なことから広く用いられている。他の手法は, データのヒストグラムを構成し, 対話的に頻度による閾値処理を繰り返して, クラスタ分割を行なうもので, カナダの CCRS で開発された (文献 2)。

本論文では, 後者と同様にヒストグラムを用いるが, その基本的考え方は, Wishart の階層的モード法 (文献 13) に近い新しい手法 (ヒストグラム・モード法と名付ける) について報告する。その基本的な考えは, 衛星画像から多次元ヒストグラムを構成し, そのモード (即ち, 確率密度の山) を多次元空間でのヒストグラム頻度のパターンとして捉え, このモードによってクラスタ解析するものである。

ヒストグラム・モード法の基礎となった階層的モード法では, 各データの近傍での標本密度から密度の極大点を求めてクラスタを構成していた。ヒストグラム・モード法では, 予めヒストグラムを構成し,  $n$  次元空間での確率密度の極大点をヒストグラムの形から決定し, その周囲にクラスタを定める。衛星画像をこのようにしてクラスタ解析するためには,

- (1) 多次元空間でのモードの把握方法
  - (a) 頻度の極大値, 谷間決定手法
  - (b) モード (クラスタ) の表現方法
- (2) 多次元ヒストグラムの効率的構成方法
- (3) 多次元ヒストグラムの標本誤差除去

などの問題に対処していかなければならない。又, このような手法的問題とは別に, 衛星画像の大規模性から必要となる膨大な計算量をどのように効率的に処理するかという計算機の高度利用も重要な課題となる。

以下, 1 章においては, 衛星画像及びその分類解析の概略と, 本論文で提案するヒストグラム・モード法の位置付けを述べる。2 章においては前記 (2) の問題を解決する, 大規模画像データの多次元ヒストグラム構成方法について述べる。3 章では前記 (1), (3) の問題を解決するヒストグラムモード法について述べる。又, 4 章では, 本プログラムにおいて用いたベクトル型計算の高度利用法について述べる。最後に, 5 章において本論文において述べる手法により衛星画像を解析した例について述べる。

## 1. 序 論

本論に入る前に, 問題の背景として, リモートセンシングに用いられる衛星画像の分類解析手法の概略について述べる。

### 1.1 衛星画像の分類解析手法

宇宙から人工衛星で地上を分光観測して得られる情報には

- (1) 地上の各地点の分光特性
- (2) 地上の表面の幾何的形狀

がある。後者の情報は得られた画像を 2 次的に観察することによって得られるもので, 地形のリニアメント抽出などが主に行なわれているが, 本論文では取り扱わない。ここでは, 前者の分光特性に基づく分類解析の問題を取り扱う。

分光特性情報は, 地上の各地点を何バンドかに分光し各地点の離散的分光スペクトルを観測して得られる。この情報から各地点のカテゴリー (即ち, 森, 都市, など) をスペクトルパターン (色) によって判別できると期待される。例えば, 図 1-1 に示すように, 森は波長の短い方にスペクトルのピークがみられ, 乾燥地面では逆に波長の長い方に高原状のピークが見られる, というように分光スペクトルパターンが異なるのでこの違いをキーに両者を弁別する。

このような分光スペクトルの違いに基づいて, 衛星画像をいくつかのカテゴリーに分類する手法は, 人間の持つ知識をどう利用するかによって, 教師付き分類 / 教師無し分類の 2 種類に分かれる。

教師付き分類とは, 判別のための弁別関数を人間

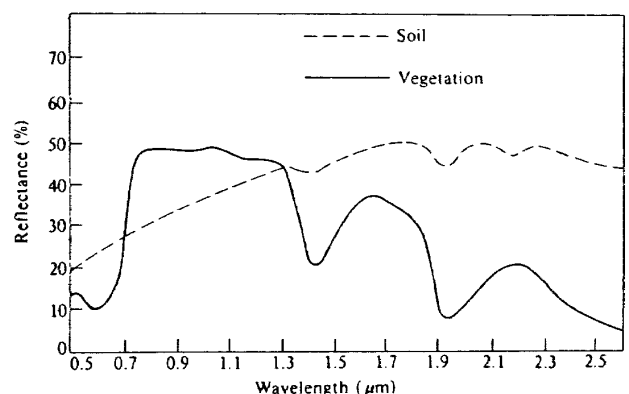


図 1-1 分光スペクトル曲線の例

(実験者, 即ち教師) の知識を直接的に利用したトレーニング過程によって定め, これによって衛星画像をカテゴリー分類するものである。この方法は簡単のため広く用いられているが, トレーニング過程の妥当性が極めて重要なファクターであるにもかかわらず, どのようなトレーニング処理が適切かについて未だ手法が確立していない難点がある。

これに対して教師無し分類は, いわゆるクラスタ解析手法を用いて, 画像からスペクトルパターンが似通っている事を基準にいくつかのクラスタを抽出し, 各クラスタを弁別カテゴリーに対応づけてカテゴリー分類するものである。実験者の知識はクラスタとカテゴリーの対応づけに用いられるだけなので, 教師付き分類に不可避であったトレーニング処理の妥当性の問題が軽減され, より客観的な解析が可能である。

クラスタ解析手法には, 従来, 解析対象のデータの質, 解析の目的などによって様々な手法が考案されている。表1-1に従来のクラスタ解析手法のうち代表的なものを衛星画像に適用した場合の問題点について示す。多くの手法では衛星画像の膨大なデータ量をうまく取り扱えず, わずかに Ball & Hall の再配置法と CCRS のヒストグラム法が実際に用いられているのみである。

再配置法を用いたプログラムでは ISODATA, ISOCLS, LARSYS 等がある。この手法では, 適

切な初期分割を与えられれば, 反復繰り返すにより最適なクラスタ分割に至る事が知られている。しかし, 適切な初期値が与えられる事が極めて重要であり, 任意の初期分割の下では準最適な解にしか収束しない事が知られている(文献6)が, 適切な初期値設定方法が未だ確立されていない点が問題である。

第2のヒストグラムによる方法は, Shlien 等によって提案されたもので, 衛星画像の多次元ヒストグラムを作り, 対話型にヒストグラムの閾値処理を行ないクラスタを分離して行くものである。この方法では処理が対話型である点に問題が残る。対話型である事は, “教師無し”である事によって得られる客観性というメリットを捨てるもので, 対話過程の妥当性という問題が新たに起きてくる。

## 1.2 ヒストグラムモード法

本論文で提案するヒストグラム・モード法は従来は衛星画像のクラスタ解析には用いられていなかった Wishart の階層的モード法を基礎にしたものである(文献13)。

階層的モード法は, それまでのクラスタ解析手法においてややあいまいなまま用いられて来たクラスタ(データの魂)の定義を, 確率密度関数のモードと定め, 理論的に明確にした。このため, 従来の手法より “自然な” クラスタを生成できる特徴を持っている。

表1-1 クラスタ解析手法を衛星画像に適用する場合の問題点

	計算負荷	メモリー量	問題点
階層的な手法	$N^2$	$N^2$	膨大な計算負荷とメモリー
非階層的な手法			
再配置法 (Ball & Hall)	$N \cdot M \cdot A$	$M$	計算負荷, 初期クラスタ
丘登り	$N \cdot M \cdot B$	$N$	計算負荷, 初期クラスタ
強制移動	$N \cdot M \cdot B$	$N$	計算負荷, 初期クラスタ
ヒストグラム法	$N' \cdot M \cdot C$	$N'$	対話型
階層的モード法	$N^2$	$N^2$	膨大な計算負荷とメモリー

$N$  : サンプル数,  $M$  : クラスタ数

$N'$  : ヒストグラムセル数 ( $10^4$ 程度)

$A, B, C$  : 計算負荷係数  $A < C \ll B$

(但し, データの再サンプリング, 次元数の減少は行なわないとした。)

階層的モード法は個体間の距離計算を基本的に用いているため、表 1-1 に示されるように計算負荷メモリーの点で少数データ用のアルゴリズムであった。

しかし、階層的モード法のアルゴリズムを詳細に検討し、

- (a) 同手法において重要な役割を果たしている“データ密度”は確率論的には確率密度関数の近似的推定である事、
- (b) クラスタの分割・融合によって探索決定される確率密度関数のモードが、多次元空間のヒストグラム頻度の極大点（即ち、密度の山）である事（図 1-2）,

の 2 点に注目し、前節のヒストグラム法と組み合わせる事によって、大規模な衛星画像にも適用できる新しい手法、ヒストグラム・モード法を以下に提案する。

ヒストグラム・モード法の基本的アプローチは、

- (1) 大規模データから多次元ヒストグラムを構成する、
- (2) ヒストグラムのモード（即ち、頻度の極大点）によって混合確率密度関数のモードを推定する、
- (3) クラスタの境界は、極大点の周囲に順次クラスタを成長させ、クラスタが互いに接触した所（即ち、頻度の谷間）とする、

に基づいている。ヒストグラム・モード法によって大規模データを効率よく解析するためには、

- 1) 多次元ヒストグラムを効率よく構成する方法

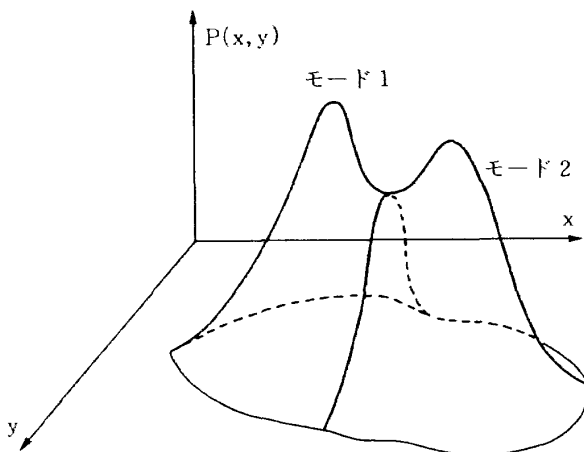


図 1-2 確率密度分布関数のモードによるクラスタ分割

（次元の呪いの問題）,

- 2) 多次元ヒストグラムの極大点を核にクラスタを構成する方法（特に計算負荷の小さい手法）,
- 3) ヒストグラムに不可避の標本誤差の処理方法, 等が問題となる。

以下、2 章において、多次元ヒストグラムの作成方法について述べ、2), 3) に関しては 3 章において述べる。

## 2. 多次元ヒストグラムの構成法 — 多層ハッシング法 —

本論文で述べるヒストグラム・モード法には多次元ヒストグラムが重要な役割を果たす。本章では多次元ヒストグラム構成に伴う問題について論じ、新たに考案した多層ハッシング法について述べる。

### 2.1 ヒストグラムの構成手法

$n$  次元ヒストグラムの構成手法には大別して 2 種類ある（図 2-1）。

- (1) 完全型  $n$  次元配列を用いる方法
- (2) 近似型
  - (a) Karhunen-Loeve 変換を用いる方法
  - (b) テーブル探索手法を用いる方法

以下、これらのそれぞれについて述べるが、その前提として、ヒストグラムを構成する対象の観測データについて次のように仮定する。

- 1) 観測データには  $N$  個の母集団  $C_i$  から生成された標本  $x$  がランダムに存在する。
- 2) 観測データは  $n$  次元ベクトル  $V$  で表わされ、各次元毎のデータは区間  $[0, a_i]$  に存在する。
- 3) 各母集団は平均  $\mu_i$ , 分散  $K_i$ , 先験確率  $P_i$  のガウス分布をなす。

このような仮定は極めて一般的であり、クラスタ解析する殆どのデータに対して仮定してよい。又、仮定 3) は 2.1.2 節において、近似的ヒストグラムに必要な領域体積を推定するためだけに必要な仮定であり、ヒストグラム構成アルゴリズムとは本質的な関係は無い。従って、2.2 節で述べる多層ハッシング法は非ガウスなデータに対しても適用できる。

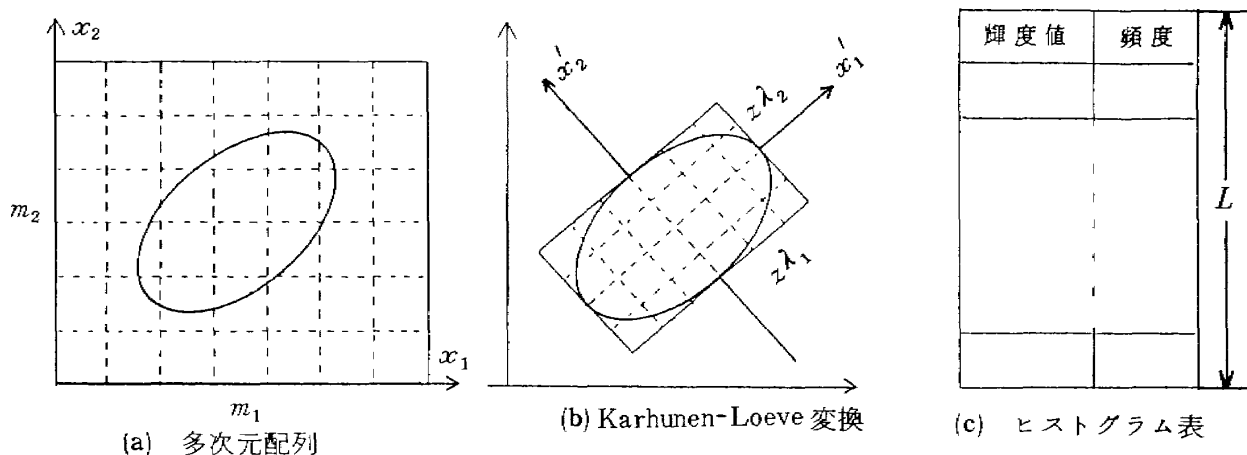


図 2-1 多次元ヒストグラムの作成法

### 2.1.1 完全なヒストグラムの構成方法

#### — $n$ 次元配列を用いる方法 —

この方法は、観測データの完全なヒストグラムを作ろうとするものである。 $n$ 次元の各データ区間  $[0, a_i]$  を適当に量子化し、離散値ベクトル  $\mathbf{v}$  を作る (式 (2.1))。離散値ベクトル  $\mathbf{v}$  の各次元毎の量子化のレベル数  $m_i$  から、ヒストグラムを構成するために用いる  $n$  次元配列は、式 (2.2) の形で定められる、いわゆる「配列」でよい。その方法に必要なメモリー量は  $V_1$  (式 (2.3)) に比例する。

$$\mathbf{v} = (v_1, v_2, \dots, v_n) \quad (2.1)$$

$$m_1 \times m_2 \times m_3 \times \dots \times m_n \quad (2.2)$$

$$V_1 = \prod_{i=1}^n m_i \quad (2.3)$$

この方法は、次元数が 1, 2 のように小さいあいだは極めて有効である。事実、2次元以下のヒストグラム作成には他の方法を考慮する必要は全くない。しかし、次元数が 3 次元以上になってくると、いわゆる「次元の呪い」のために、ヒストグラム構成に必要なメモリー領域が天文学的になりはじめ、4 次元以上ではこの方法は殆どの場合役に立たない。

データは多次元空間の全てに一樣に存在するわけではなく、局所的にかたまって存在している (従ってクラスタ解析を行なう意味がある)。このため、この莫大なメモリー領域の殆どは結果的に頻度 0 で本質的には不必要なメモリー領域である。

### 2.1.2 近似的ヒストグラムの構成

観測データの「完全な」ヒストグラムを作ろうと

すると 2.1.1 節で述べたように「次元の呪い」のために莫大なメモリー領域が無駄になるが、ヒストグラムの完全性を多少犠牲にした近似的なヒストグラムならば、莫大な無駄領域のほとんどを必要としない手法が開発できる。

観測データになんらかの統計的母集団を想定すれば、観測データは多次元空間のある部分に密に集中しているはずで、他の大半の領域では極めて疎にしか存在せず無視し得るであろう。従って、データが集中する部分でだけ集中的にヒストグラムを作り、他の部分は無視しても観測データの大多数を表わすヒストグラムが構成でき、これによって全空間のヒストグラムが近似できるであろう。

各母集団からの標本がその先験確率に比例してヒストグラムに含まれるものとし、ガウス分布を仮定すると、各母集団の平均の周りに部分領域  $A_i$  が次式で定められる。

$$A_i = \{X \mid C_i \ni X,$$

$$(X - \mu_i)^t K_i^{-1} (X - \mu_i) < \chi^2(n, 1-P)\} \quad (2.4)$$

$(X - \mu_i)^t K_i^{-1} (X - \mu_i)$  は自由度  $n$  の  $\chi^2$  分布をなすので、式 (2.4) で定まる領域  $A_i$  には母集団  $C_i$  から生成される標本  $X$  のうち  $P\%$  のものが含まれる。従って、これらの領域  $A_i$  ( $i = 1, \dots, N$ ) にだけしかヒストグラムセルを設定しなくても、確率的に全データの  $P\%$  に相当するデータからなるヒストグラムが構成できる。必要なヒストグラムセルの数は領域同士の重なりを無視すれば、各領域  $A_i$  の体積の総

和  $V_2$  に比例する。

$$W_i = \frac{(\pi z^2)^{n/2} \cdot |K_i|^{1/2}}{\Gamma(n/2 + 1)}, i=1 \cdots N. \quad (2.5)$$

$$V_2 = \sum_{i=1}^N W_i = \frac{(\pi z^2)^{n/2}}{\Gamma(n/2 + 1)} \cdot \sum_{i=1}^N |K_i|^{1/2} \quad (2.6)$$

$$\text{但し } z^2 = \chi^2(n, 1-P)$$

$$R_{12} = \frac{V_1}{V_2} = \prod_{i=1}^n \frac{m_i}{(\pi z^2)^{1/2}} \cdot \frac{\Gamma(n/2 + 1)}{\sum_{j=1}^N |K_j|^{1/2}} \quad (2.7)$$

完全なヒストグラムを作る場合に必要なた体積との比は図2-2に示すように、次元数の冪で大きくなる。従って、近似的ヒストグラム構成法は完全なヒストグラムを作る場合よりもメモリー領域の有効利用が

図れ、より高次元なデータに対して適用できると言える。

近似的ヒストグラムを構成する方法として、

Karhunen-Loeve 変換を用いる方法

テーブル探索手法を用いる方法

が考えられる。本論文で述べる多層ハッシング法は後者のテーブル探索による手法のうちハッシング法を発展させたものである。

#### (1) Karhunen-Loeve 変換を用いる方法

この方法は2.1.1節の多次元配列法を用いて近似的ヒストグラムを構成する方法である。

ヒストグラムを作るときに各母集団の統計量が既知であるということは、一般には仮定できない。従って、領域  $A_i$  を個別に予め定めることはできない。しかし全観測データの平均  $\mu_0$ 、共分散  $K_0$  は簡単に求められるから、各個別領域  $A_i$  の代わりに全体で

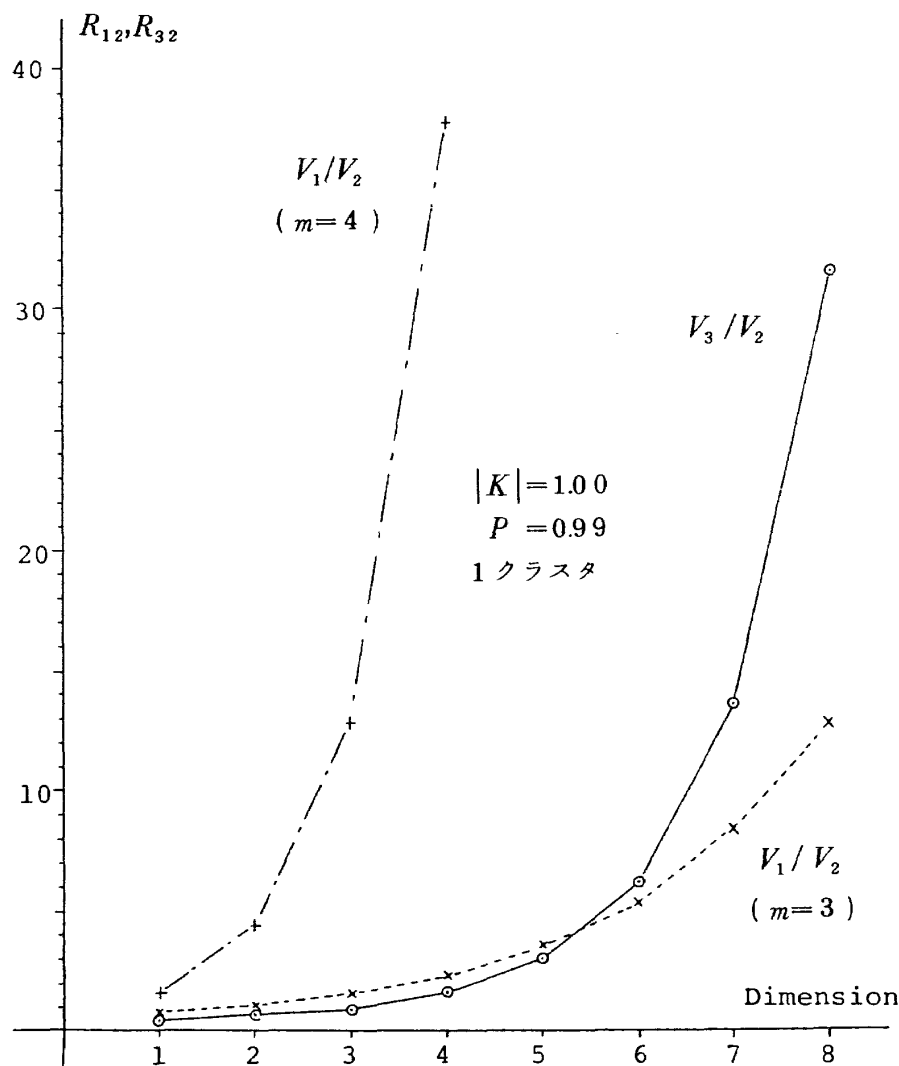


図2-2 多次元ヒストグラムを構成するのに必要なメモリー量の比較

1つの領域 $A_0$ を $\mu_0, K_0$ によって定められる。このときヒストグラムの近似度 $P$ がある程度大きければ $A_0$ によって各個別領域 $A_i$ の和領域 $A_T$ を近似してもよからう(付図1参照)。

$$A_T = \bigcup_{i=1}^N A_i \quad (2.8)$$

$$V_3' = \frac{(\pi z^2)^{n/2}}{\Gamma(n/2+1)} \cdot |K_0|^{1/2} \quad (2.9)$$

$A_0$ は $n$ 次元楕円体であるが、Karhunen-Loeve変換には、これに外接する $n$ 次元直方体 $A_0'$ を用いる。 $A_0'$ はその体積が $V_3'$ (式(2.10))で表わされ $A_0$ と殆ど同程度であり、かつアドレス計算は極めて簡単である。

$$V_3' = (2z)^n \cdot |K_0|^{1/2} \quad (2.10)$$

$$R_{33} = \frac{V_3'}{V_3} = \frac{2^n \cdot \Gamma(n/2+1)}{(\pi)^{n/2}} \quad (2.11)$$

観測値 $X$ を、 $\mu_0, K_0$ によってKarhunen-Loeve変換し $X'$ を得る(式(2.12)) $T$ はKarhunen-Loeve変換行列。 $n$ 次元直方体領域 $A_0'$ は、共分散行列 $K_0$ の固有値 $\lambda_i$ によって式(2.13)で規定されるから、各軸毎に区間 $[-Z\lambda_i, Z\lambda_i]$ を所要の分解能に従って適当に量子化してヒストグラムを作成すればよい(図2-1(b))。

$$X' = T \cdot (X - \mu_0) \quad (2.12)$$

$$A_0' = \{X' \mid -Z\lambda_i < x'_i < Z\lambda_i\} \quad (2.13)$$

この方法に必要なメモリー領域の大きさは式(2.10)の $V_3'$ に比例するが、これは $V_2$ に比べると(式(2.14))、母集団の数 $N$ 及び $\mu_i, K_i$ によって変わるが、数倍程度は大きいと言わざるを得ない。

$$R_{32} = \frac{V_3'}{V_2} = \frac{2^n \cdot \Gamma(n/2+1)}{(\pi)^{n/2}} \cdot \frac{|K_0|^{1/2}}{\sum_{i=1}^N |K_i|^{1/2}} \quad (2.14)$$

## (2) テーブル探索手法を用いる方法

本節で述べる手法は、2.2節で述べる多層ハッ

シング法の基礎をなす手法である。

Karhunen-Loeve変換を用いてヒストグラム構成範囲をデータが密に存在する領域に限ってメモリー領域の節約を図る前者の方法は、力まかせの多次元配列法に比べればはるかにメモリーを効率よく使えるが、それでも $P\%$ の観測データに相当するヒストグラムを作るのに本当に必要な領域の広さ $V_2$ と比べると無駄な領域が残る。

本節では原理的にメモリー領域の広さが $V_2$ に比例するテーブル探索型のヒストグラム構成方法について述べる。この方法はヒストグラム表及びテーブル探索手法からなる。ヒストグラム表(図2-1(c))は、

$$(\text{頻度値}, \text{観測データ}) = (Hf(i), Hv(i))$$

を1セットにしたヒストグラムセルを1つの要素とする長さ $L$ の表である。表の長さ $L$ は、データ $v$ の取り得るあらゆる場合の数 $V_1$ よりずっと小さく定める。表の各要素とデータ $v$ の間には1:mの写像関係が成り立つ。このため、あるヒストグラムセルの位置 $k$ に対応する観測データ $v$ は何通りも存在し、ヒストグラムセルの衝突が起きる。この衝突の問題を解決し、観測データ $v$ からセルの位置 $k$ を決定するために、テーブル探索手法が用いられる。

ハッシング法は、このような $m:1$ の写像関係を効率よく取り扱えるテーブル探索手法として知られている。ハッシング法では、ヒストグラム表の長さ $L$ を素数 $P_1$ になるように選ぶ。 $P_1$ に近く $P_1$ を超えない他の素数を $P_2$ とする。観測データ $v$ は式(2.1)のように量子化されているとする。この時、ハッシング法では、ヒストグラム表探索のためのキー $a_i, b$ を式(2.16)で定め、これを用いてStep 1-Step 3(図2-3(b))のアルゴリズムで観測データ $v$ を格納するヒストグラムセルを探索する。

$$V = (\cdots (v_1 * m_1 + v_2) * m_2 + v_3) * m_3 + \cdots * m_n + v_n \quad (2.15)$$

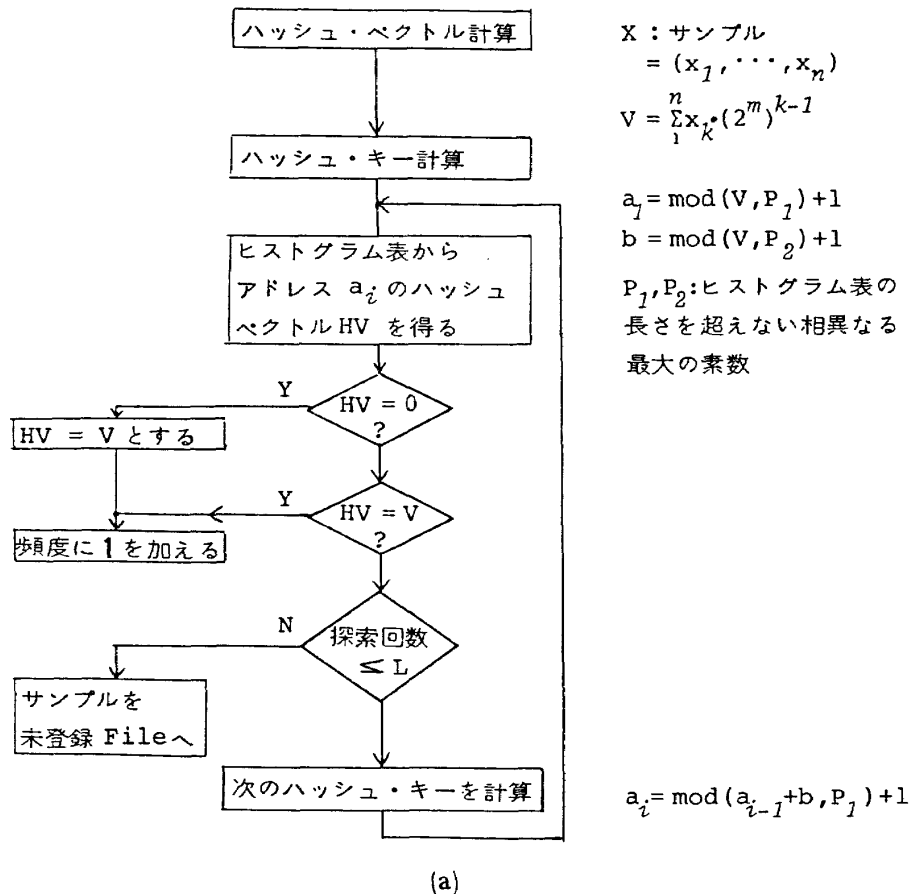
$$a_1 = \text{mod}(V, P_1) + 1$$

$$b = \text{mod}(V, P_2) + 1 \quad (2.16)$$

$$a_i = \text{mod}(a_{i-1} + b, P_1) + 1, i \geq 2$$

図2-3にハッシング法のアルゴリズムを示す





(a)

Step1: 観測データ  $v$  からハッシュ・キー  $a$ ,  $b$  を求める。

Step2: ヒストグラム表の  $a$  番目の要素を調べ、

- $H_f(a) = 0$  ならば  
 $H_v(a) = v$ ,  $H_f(a) = 1$  として,  $v$  の処理終了。
- $H_v(a) = v$  ならば  
 $H_f(a) = H_f(a) + 1$  として,  $v$  の処理終了。
- それ以外ならば, Step 3 へ。

Step3:  $i = i + 1$  として、

$i$  が  $N$  以下ならば, 式 (2.16) で  $a_i$  を求め

Step 2 へ、

$i$  が  $N$  以上ならば,  $v$  を未登録データとして終了。

(b)

図 2-3 ハッシング法によるヒストグラム作成のアルゴリズム

(これを以下基礎的ハッシング法と呼ぶ)。

Shlien 等はこのハッシング法をランドサット M SS 画像に応用し、ヒストグラム表の長さ 6,000 で 250,000 画素の 95% 以上を表わす 4 次元ヒストグラムが殆どのシーンに対して、構成できる事を示した。

## 2.2 多層ハッシング法

本節では多次元ヒストグラムを効率よく作成するために、筆者等が新たに考案した多層ハッシング法について述べる。

### 2.2.1 多層ハッシング法

ハッシング法は、本来はハッシング表に登録される相異なる事象の数が、表の長さ  $L$  より小さい場合

のアルゴリズムである。このためハッシング表（ここではヒストグラム表）が満杯になった時の処理アルゴリズムは不十分なものしかない。しかし、衛星画像のような膨大なデータの多次元ヒストグラムを、比較的小規模なヒストグラム表によって作る場合には、ハッシング表が満杯になる事は不可避免なのでその場合を効率的に取り扱う手法がハッシング法に付加されなければならない。

従来はこの問題は殆ど論じられておらず、わずかに

「表が一杯になりかけたら、頻度の少ないヒストグラムセルを取り除いて空き領域を作り、残ったヒストグラムセルを再度ハッシングしてヒストグラム表を作り直す。」（再ハッシング処理）

と言う簡単な方法が知られているのみであった。このような方法によってもデータ量が少ない間は余り重大な問題にはならなかった。しかし、データ量が増して来た時

- (1) ある程度以上のデータを処理した後は、頻度の少ないヒストグラムセルを除く再ハッシング処理が頻般に必要なになる。この結果、計算時間が急激に増大し、ヒストグラムの計算負荷よりも再ハッシング処理の計算負荷の方がはるかに大きくなる。
- (2) ヒストグラムセルを取り除くときに、これから処理するデータにおける頻度分布が全く考慮されない。このため、実際には頻度の大きなセルであるにもかかわらず、構成されたヒストグラムに存在しない事が在り得る。

と言う問題がデータ量が増すにつれ重要になって来る。特に、第2の点は衛星画像においては重要である。各クラスタ（母集団）のデータは衛星画像の中

で空間的に偏って存在しているが、ハッシング法では結果が処理にける画像データの順序に依存するので、画像の後半部分に存在する大きなクラスを逃す事につながってくる。

### 2.2.2 多層ハッシング法のアルゴリズム

ここで述べる多層ハッシング法は、このような問題を多数のヒストグラム表を動的に使用する事によって解決するものである。

多層ハッシング法では図2-4に示すように、最初のステップでは衛星画像の帯状の部分画像のそれぞれに対して別々のヒストグラム表を用いてヒストグラムを作る。部分画像の大きさは一定ではなく、画素当りの計算時間の増加、ヒストグラム表への未登録画像数の増加（図2-3 Step3）等、ハッシングの効率が低下して来た場合に、新しいヒストグラム表に切り替える。このようにして衛星画像全体に対応するヒストグラム表を作った後に、次のステップでこれらのヒストグラム表を融合して1つのヒストグラムとする。ヒストグラム表の融合にもやはりハッシング法を用いる。

基礎的ハッシング法で問題となった再ハッシング処理は多層ハッシング法では行わず、複数のヒストグラム表の使用とそれらの融合によって計算負荷を軽減する事が出来る。

又、クラスタが空間的に偏在しても、多層ハッシング法では最初のステップで部分画像毎にヒストグラム表を更新するので、クラスタが十分な割合を占めていさえすれば、いずれかのヒストグラム表で大きな頻度を持つヒストグラムセルとなり、最終的なヒストグラムに適切に反映される。

多層ハッシング法のアルゴリズムを以下に示す

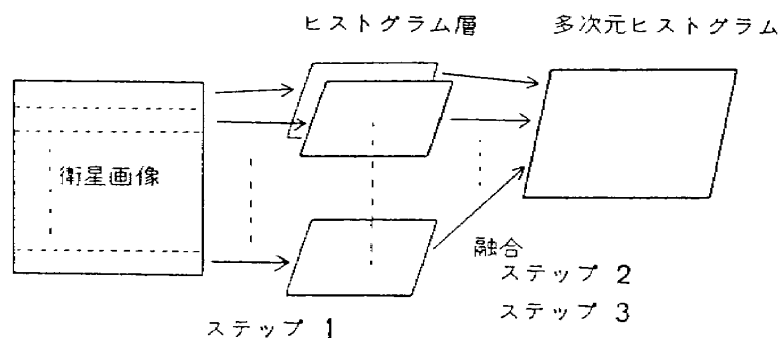


図2-4 多層ハッシング法

(図2-5, 6, 7)。本手法の基本的考えは、大量の生データ(衛星画像)を取り扱う部分(Step A)は、できるだけ効率の良い高速処理を行ない、ある程度圧縮したデータ(ヒストグラムセル)に対しては手が込んで時間がかかるが、高精度な処理

(Step B, C)を行なうものである。このような高速処理と高精度処理の組み合わせにより、処理時間の短縮が図れるばかりでなく、衛星画像に見られるクラスターの空間的な偏りにも効果的に対処できる。

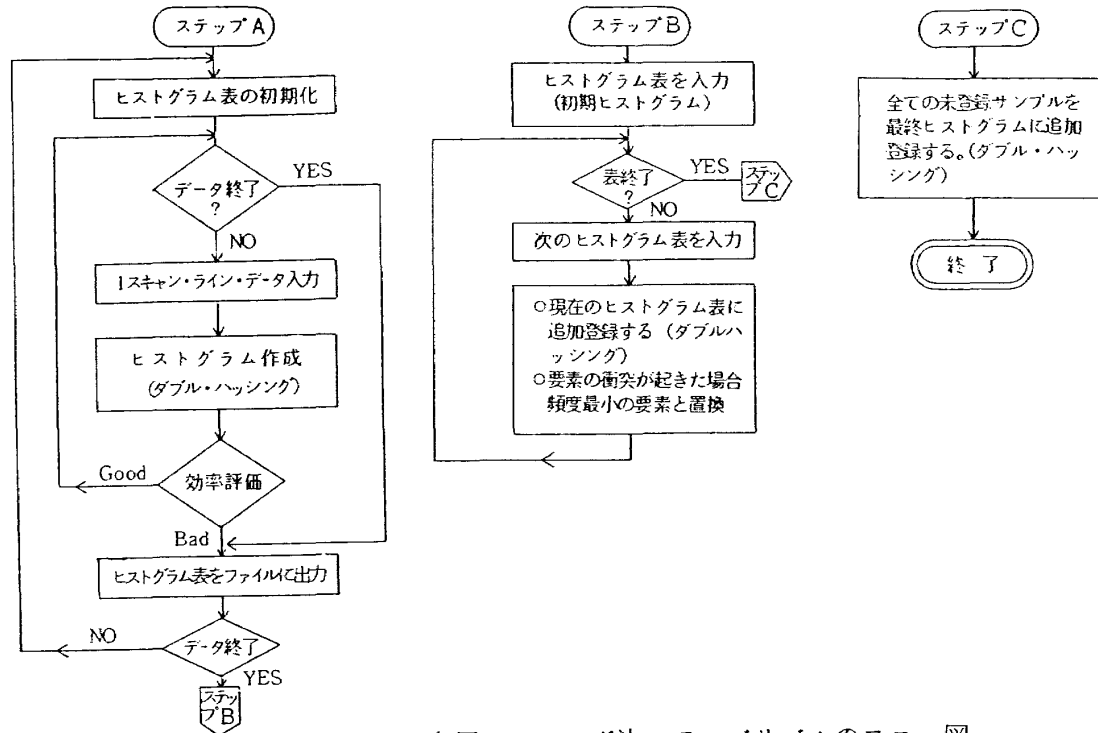


図2-5 多層ハッシング法 アルゴリズムのフロー図

#### [ Step A ]

Step A-1: 衛星画像の1スキャンラインのデータを取り、従来のハッシング法で多次元ヒストグラムを作る。ハッシングで未登録となったデータはFile Bにしまう。

Step A-2: ハッシング法の効率を計算し、一定の条件が満たされた時、Step 3に行く。データが終了した時はStep Bへ行く。それ以外はStep 1に行く。

Step A-3: 現在のヒストグラム層を終結しヒストグラム表をFile Aにしまう。新しいヒストグラム表を次のヒストグラム層作成のために用意する。

#### [ Step B ]

Step B-1: File A からヒストグラム表を1つ取り出し、現在のヒストグラム表とする。

Step B-2: File A から次のヒストグラム表を取り出し、現在のヒストグラム表にハッシング法に従って登録して行く。“衝突”の場合は衝突の連鎖上でセルの頻度を比較し頻度最小のセルを現在のヒストグラム表から取り除き、File Bにしまう。

Step B-3: File A にヒストグラム表がなくなるまでStep B-2を繰り返す。

#### [ Step C ]

Step B で構成したヒストグラム表にFile Bにしまわれている画素及びヒストグラムセルをハッシング法で登録する。

図2-6 多層ハッシング法のアルゴリズム

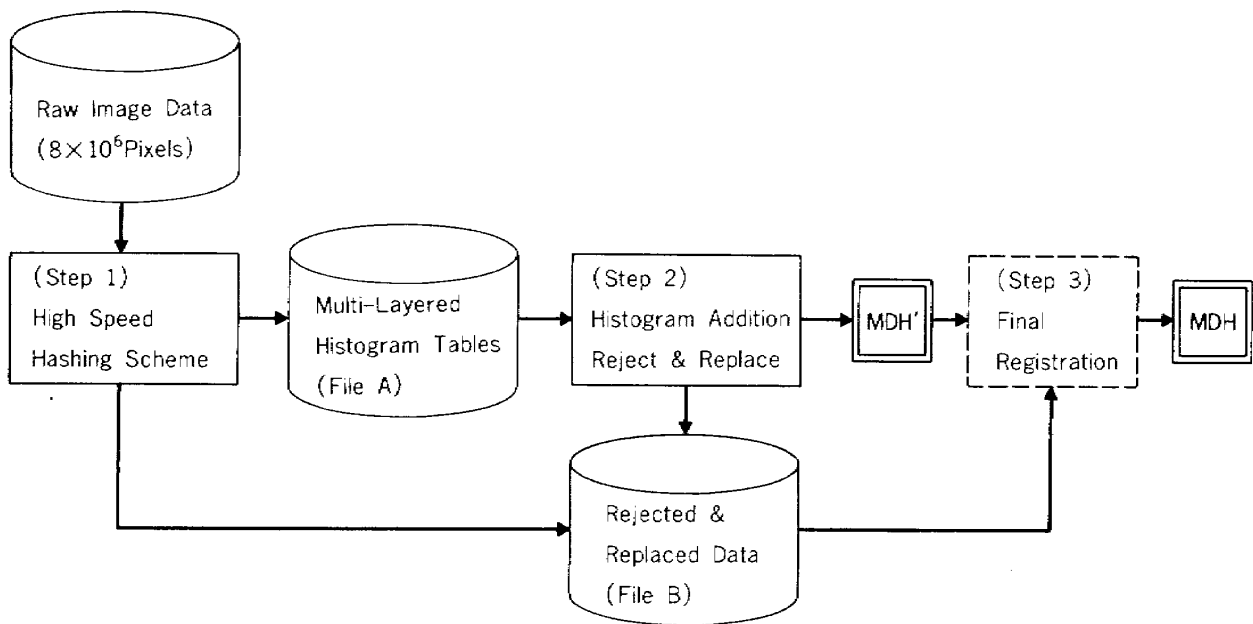


図2-7 多層ハッシング法

### 3. ヒストグラム・モード法

#### 3.1 ヒストグラムモード法

本論文で述べるヒストグラム・モード法はWishart の階層的モード法の考え方を基礎にしたものである。階層的モード法は計算負荷、メモリーの点で小規模データにしか適用できなかったが、ヒストグラム・モード法では、多次元ヒストグラムを用いる事によって必要な計算負荷、メモリー量が大幅に軽減され、これにより大規模な多次元衛星画像を一度にクラスタ解析できるようになった。

##### 3.1.1 階層的モード法の大規模化の検討

階層的モード法のアルゴリズムを図3-1に示す。階層的モード法では確率的分布状態を把握するために、各点の近傍でのデータ密度を

“ $k$ 番目に近い点までの距離”

によって定義する (Step-1)。ここで $k$ は適当に定めるパラメータで、通常2~5が用いられる。距離のリスト $PD$ を作るためには、データ数を $N$ として、 $N(N-1)$ 回の距離計算及び比較計算が必要になる。又、Step-2の並べ換え計算は本質的に不可避であり、このためには $N \cdot \log N$ のオーダーの計算が必要である。更に、Step-4のクラスタを構成していく処理でも $N(N-1)$ 回の距離計算及び

比較計算が原理的に必要である。

このように、階層的モード法の基本的アルゴリズムは計算負荷が標本数 $N$ の自乗に比例するという点で、大規模データのクラスタ解析には向いておらず、主として少数標本の場合に向けて考えられていると言えるであろう。

しかし、階層的モード法のアルゴリズムを更に詳しく検討した結果、多次元ヒストグラムを用いる事によって、大規模データにおいても同様の手法が導出できる事が以下の如く明らかになった。

(a) この手法において重要な役割を果たす“データ密度”とは、確率論的には確率密度関数の値と同義であるから、各点のデータ密度の指標配列 $PD$ の代わりに各点近傍でのヒストグラムの頻度値を用いても同じ結果が得られる。

(b) 多次元空間でのヒストグラムの構成方法については第2章で詳しく述べたが、その計算負荷は $N^2$ ではなく $N$ のオーダーである。又、ヒストグラムを構成する際に2.2節で述べた多層ハッシング法を用いれば、 $10^6$ のオーダーの標本を $10^4$ 個程度のヒストグラムセルで表わせる。このため、Step-2のソーティングの計算負荷( $N \log N$ )も大幅に軽減できる。

(c) 最後に、クラスタ分割生成のStep-4-a~cのアルゴリズムは、基本的には、クラスタの核を

Step-1: データ間の距離を計算し、各点 $i$ 毎に各点から $k$ 番目に近い点の番号 $\alpha_i$ とその距離 $\beta_i$ のリスト $PD(i) = \beta_i$ , $KP(i) = \alpha_i$ を作る。	「密な点リスト」上の各点との距離を計算し、
Step-2: 距離リスト $PD$ を距離の小さい順にソートする。点リスト $KP$ も $PD$ と同じ順に並べかえる。	(a) 全ての距離が $PMIN$ 以上なら点 $a$ によって新しいクラスタを作る。
Step-3: $KP(1)$ を第 1 のクラスタとする。クラスタ 1 の密な点のリストに $KP(1)$ を加える。 $j = 2$ とする。	(b) 1 つのクラスタの密な点からだけ $PMIN$ 以内の距離にあるならば、そのクラスタの点として「密な点リスト」に付加する。
Step-4: $KP(i)$ の点を $a$ とし、 $PMIN = PD(j)$ とする。点 $a$ とそれまでに作られたクラスタの	(c) 2 つ以上のクラスタの密な点から $PMIN$ 以内の距離にあるならば、それらのクラスタと点 $a$ が 1 つのクラスタに融合される。
	Step-5: $j = j + 1$ として $KP$ が尽るまで Step-4 を繰り返す。

図 3-1 階層的モード法のアルゴリズム

“データ密度の高い点”に定め、データ密度がそれよりも低い点を順次その核の周りに付加していく (Step-4-b) 事と、新たなクラスタの設定を、新たに付加する点がそれまでのクラスタから遠く離れている場合 (Step-4-a) に行なう事からなっている。このアルゴリズムは多次元空間において、“データ密度 (= ヒストグラム頻度)” の極大点を抽出するアルゴリズムであるから、より計算負荷の小さい極大点探索手法に置き換えてもよい。

### 3.1.2 ヒストグラム・モード法

本手法では、前節で述べた考えに基づき、多次元ヒストグラムの形 (極大点) によってクラスタを決定する。

大量のデータからヒストグラムを作った場合、ヒストグラムによってそのデータの母集団の確率密度関数を推定する事が出来る。これは、元のデータが複数の母集団からの標本が混合したものであっても同様で、この場合は混合確率密度関数がヒストグラムから推定される。各母集団の確率密度関数が単峰であると仮定すると、母集団同士の分布関数の重なり程度が低いならば、これら  $M$  個の母集団によって作られる混合確率密度関数には  $M$  個の極大点が現われ、Bayes の弁別境界はこれら  $M$  個の極大点の間の確率密度の谷間となる。(図 1-2) (厳密には次節に述べるように若干異なる)

従って、ヒストグラムの極大及び谷間の位置によ

って混合確率密度関数の極大・谷間位置を推定し、これによってクラスタ分割をすれば、ほぼ最適なクラスタ分割が得られる、と言えるであろう。

ヒストグラム・モード法はこのような考えに基づき、衛星画像データの多次元ヒストグラムの極大点と谷間を探索決定する事によって、クラスタ分割を行なうものである。

### 3.1.3. ヒストグラムの谷間と最適弁別境界

ヒストグラム・モード法ではヒストグラム頻度の谷間によってクラスタを分割する。しかし、ヒストグラムの谷間、即ち確率密度の谷間は、厳密には Bayes の最適弁別境界にならないので、その違いを以下に例をとって示す。

簡単のために、データを 1 次元、2 母集団とする。各々の母集団  $A$ ,  $B$  の確率密度関数  $P_a(x)$ ,  $P_b(x)$  が、

- (1) 三角分布
- (2) ガウス分布

の 2 通りについて示す。

- (1) 三角分布の場合:

確率密度関数  $P_a(x)$ ,  $P_b(x)$  は

$$P_a(x) = \begin{cases} 0 & \cdots a_1 < x \\ \frac{P_a}{a_1^2} (a_1 - x) & \cdots 0 < x \leq a_1 \\ \frac{P_a}{a_1^2} (a_1 + x) & \cdots a_1 < x \leq 0 \\ 0 & \cdots x \leq -a_1 \end{cases}$$

$\alpha = 1$  として

$$b_0 = \frac{a_1 a_3}{a_2 + a_1} \quad (3.5)$$

これに対して、ヒストグラムに現われる頻度の谷間点  $b_1$  は図 3-2 では

$$b_1 = a_1 \quad (3.6)$$

この差  $\Delta$  は

$$P_b(x) = \begin{cases} 0 & \cdots a_2 + a_3 \leq x \\ \frac{P_b}{a_2^2} (a_2 + a_3 - x) & \cdots a_3 < x \leq a_2 + a_3 \\ \frac{P_b}{a_2^2} (a_2 + a_3 + x) & \cdots a_3 - a_2 < x < a_3 \\ 0 & \cdots x \leq a_3 - a_2 \end{cases} \quad (3.1)$$

$$\Delta = \frac{(b_0 - b_1)}{a_3} = \frac{a_1}{a_3 (a_1 + a_2)} (a_3 - (a_1 + a_2)) \quad (3.7)$$

となる。従って分布  $A, B$  の重なりが少なければ、即ち  $a_3 - (a_1 + a_2)$  が小さければ、最適弁別境界  $b_0$  と頻度の谷間点  $b_1$  の差  $\Delta$  は小さく無視し得る。

(2) 正規分布の場合：

2つの母集団の確率密度関数を

とする (図 3-2)。その混合確率密度関数  $P_0$  は

$$P_0(x) = P_a(x) + P_b(x) \quad (3.2)$$

であり、図 3-2 の点線のような分布になる。

この時、Bayes の最適弁別境界点  $b_0$  は、クラス  $A, B$  の損失係数を等しいとして

$$\int_{b_0}^{a_1} P_a(x) dx = \int_{a_3 - a_2}^{b_0} P_b(x) dx \quad (3.3)$$

から求められ、

$$b_0 = \frac{a_1 a_2 (1 - \alpha) + a_1 a_3 \alpha}{a_2 + a_1 \alpha}$$

$$\text{但し } \alpha = (P_b / P_a)^{1/2} \quad (3.4)$$

$$\begin{aligned} P_a(x) &= \frac{P_a}{(2\pi)^{1/2} \sigma_1} \cdot \exp\left(-\frac{x^2}{2\sigma_1^2}\right) \\ P_b(x) &= \frac{P_b}{(2\pi)^{1/2} \sigma_2} \cdot \exp\left(-\frac{(x-d)^2}{2\sigma_2^2}\right) \\ P_a &= P_b, \quad \sigma_1 = \sigma_2 \end{aligned} \quad (3.8)$$

の時以外は、最適弁別境界  $b_0$  と頻度の谷間点  $b_1$  は若干異なる。

図 3-3 に一次元の場合のいくつかの計算例を示す。

### 3.1.4 クラスタの定義

リモートセンシングにおけるクラスタ解析の目的は、観測データからある特徴をもつ代表的クラスタ

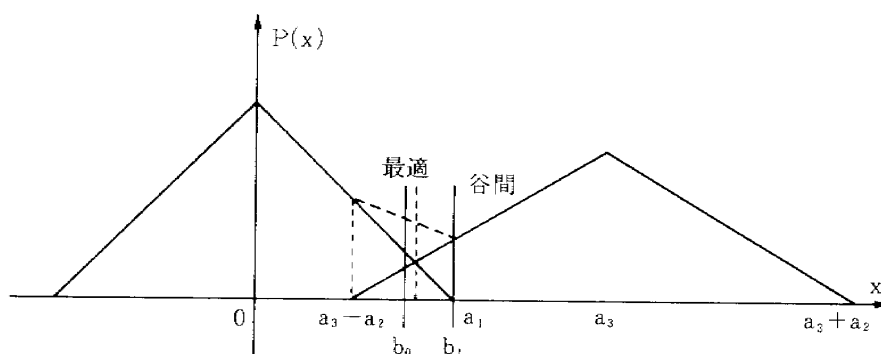


図 3-2 一次元三角分布の弁別境界

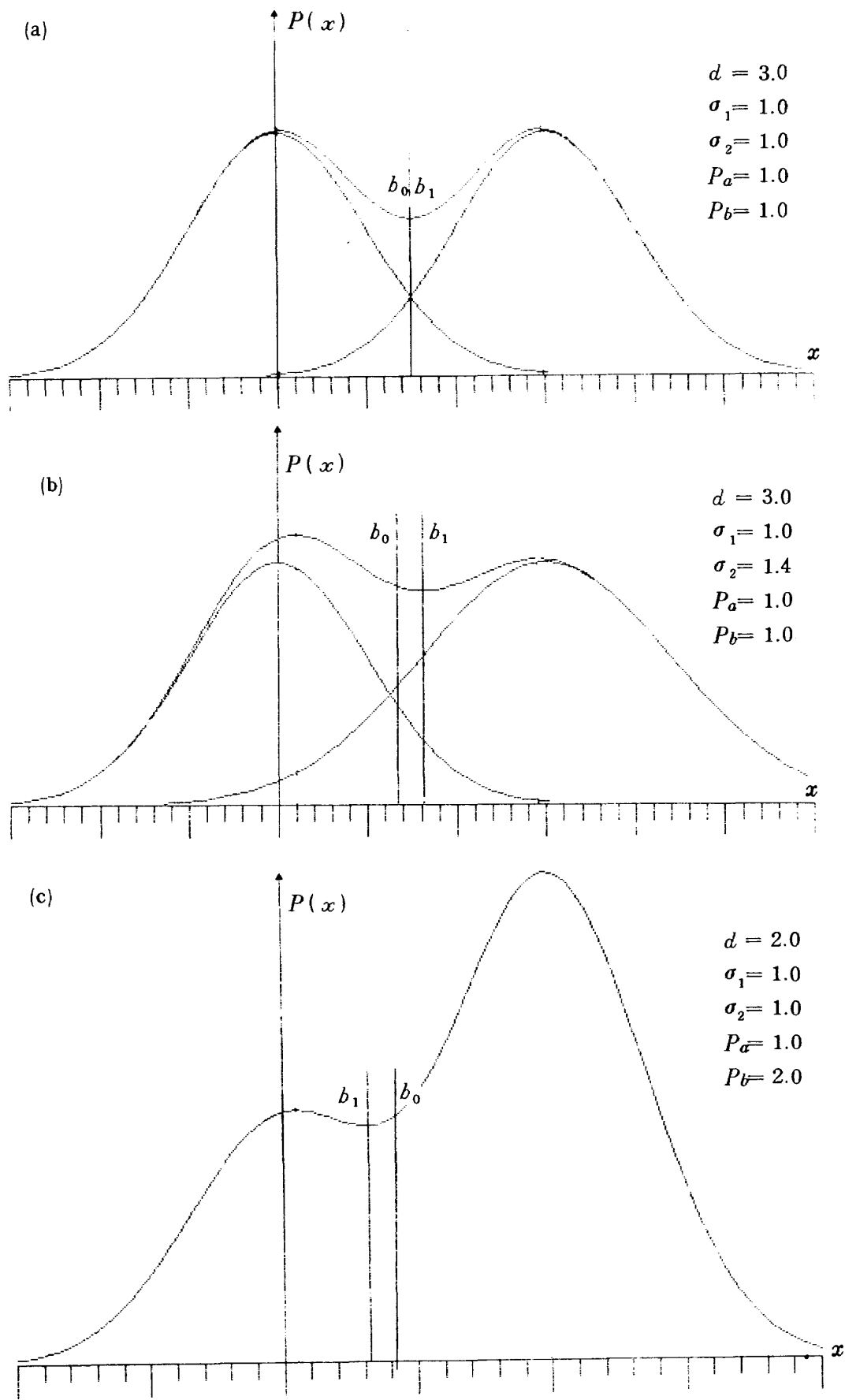


図 3-3 一次元正規分布の弁別境界

を抽出し、これによって観測画像の弁別分類を行なう事であった。抽出されたクラスタが目的（リモートセンシングデータ解析）にあっていないか否かは、最終的には弁別分類の結果であるクラスタ図を実験者が判定して決める。

しかし、クラスタ解析の処理過程ではこの判定は行なえないので、クラスタ解析処理が正しく行なわれているか、又、いつ処理を終了するかなどの判定にクラスタの妥当性の基準が必要である。更に、妥当性の基準が一般性にとみ、合理的なものであればあるほど、最終的なクラスタ図の妥当性が強くなり、解析処理の独自性が発揮される。

クラスタの妥当性基準としては、従来

- (a) minimum variance 基準
- (b) エントロピー基準

等の、全体的な傾向を示すなんらかの統計量を求め、それによって定量的に評価されていた。しかし、このような定量的評価によって得られるクラスタが統計論的にどのような意味を持つのか若干あいまいであった。

Wishartの階層的モード法は、この点を明確にするもので、クラスタ解析によって確率密度関数のモードを把握する事を目的にアルゴリズムが構成されている。

確率密度関数のモードとは図1-2に見られるような、確率密度の極大点を中心とし、確率密度の谷間を隣接する領域との境界とする領域である。このようなクラスタの定義は、統計論的に明確であるばかりでなく、定性的にも極めて受け入れやすいものである。

本論文で述べる、ヒストグラム・モード法においてはクラスタは、

ヒストグラムの極大点を核とし、

頻度の谷間を隣接するクラスタとの境界とする  
多次元空間の領域、

と定義する。これは、確率密度関数のモードの概念を含み、さらに複雑な形状の領域をも含むように拡張された定義である。クラスタは単峰でありさえすれば十分で、正規分布以外にも三角分布や馬蹄形のクラスタもこの定義に含まれ、ヒストグラム・モード法によって抽出され得る。

この事はリモートセンシングにおいては特に重要なメリットである。衛星画像から抽出されるクラスタが、必ずしも正規分布していない可能性が指摘されているので（文献5）、単純な形以外のクラスタもそのまま抽出できる事は重要である。

### 3.1.5 ヒストグラム・モード法の限界

ヒストグラム・モード法では、クラスタ同士をヒストグラム頻度の谷間で分割する。このことは、仮定：N個の単峰なクラスタが観測データに存在するならばその混合確率密度関数はN峰となる、ということが大前提として仮定されていることを意味する。

この仮定が成立しない場合には、ヒストグラム・モード法によってはクラスタの分割ができない。このような場合の若干の例を図3-4に示す。

図3-4に見られるように、クラスタ同士が近すぎ過ぎ混合確率密度関数の形が高原状になった場合には、ヒストグラム・モード法ではこれを2つのクラスタに分割できない。しかし、この場合でも、第3のクラスタとの識別は行なえる。従って、クラスタ分割をしたのち、必要ならばクラスタ個々の妥当性を正規性検定等の統計的検定によって検証すればよからう。

## 3.2 ヒストグラム・モード法のアルゴリズム

### 3.2.1 概念の定義

ヒストグラム・モード法で用いる概念を定義する。

- (1) クラスタ：“クラスタ候補”のうち“クラスタ確立条件”を満たしたものを“クラスタ”と呼ぶ。クラスタは他のクラスタと融合されない。
- (2) クラスタ確立条件：“クラスタ候補”を取り巻く谷間の深さの最小値が $N_v$ 以上である事。 $N_v$ は谷間の深さの閾値。
- (3) クラスタ候補：“孤立点”の周囲に順次拡張融合された“領域”の集合。
- (4) 孤立点：多次元空間の一点で、その点が処理されるまでに構成された、どの“領域”とも“連結”しない点。孤立点の周りには新しいクラスタ候補の核となる領域が作られ、その領域だけからなるクラスタ候補となる。



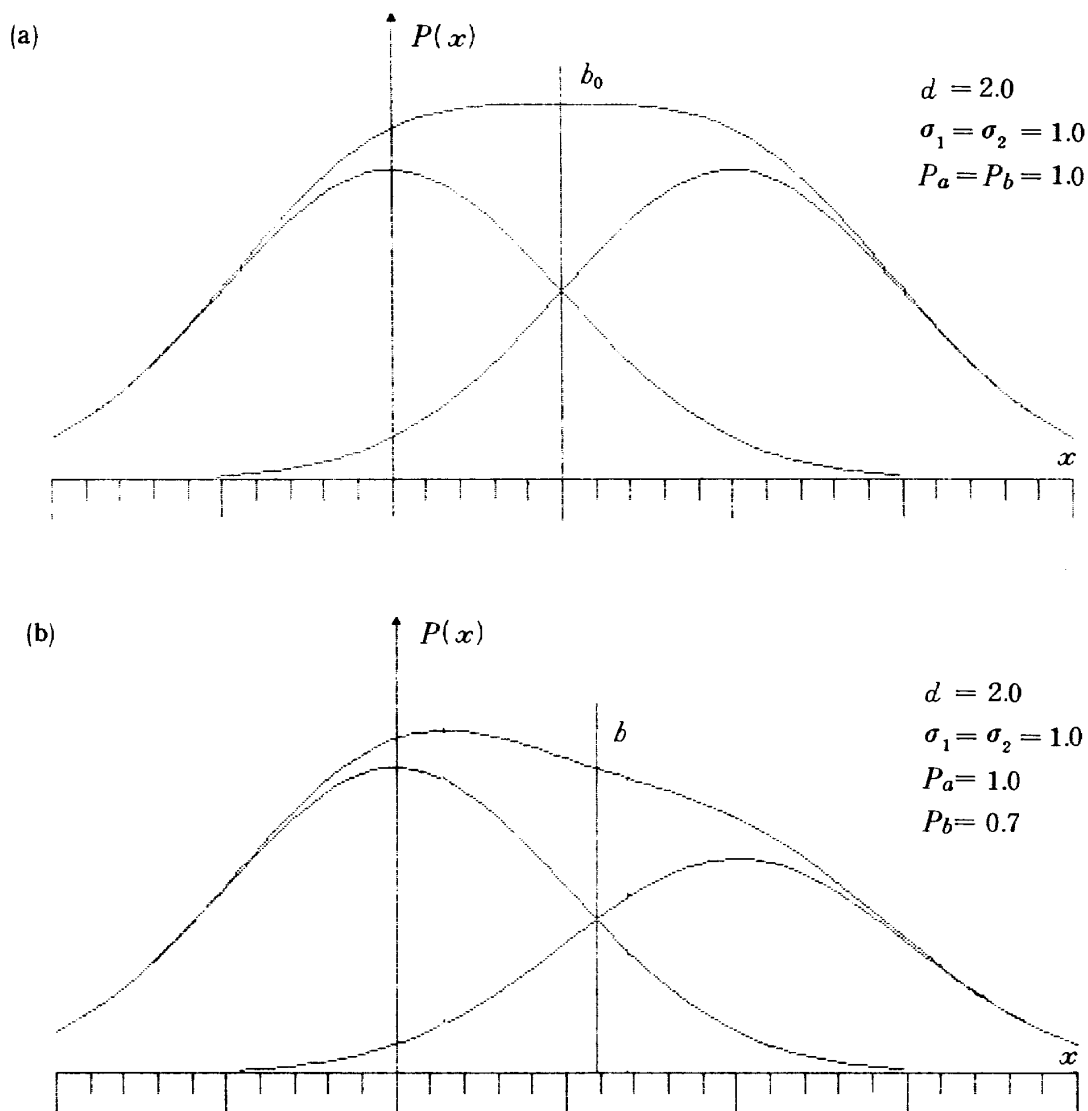


図 3-4 ヒストグラム・モード法によって分割できないクラスタ例

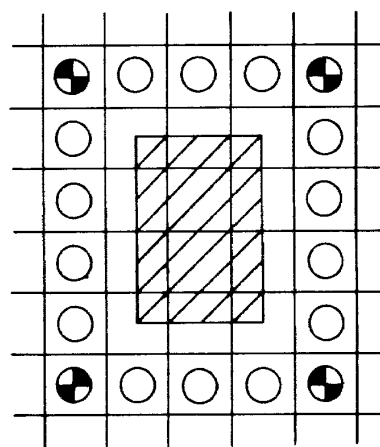
- (5) 領域：クラスタ/クラスタ候補の要素。領域に含まれるヒストグラムセルの集合に外接する，座標軸に平行な多次元直方体であり，各次元毎の上・下限  $R_i^v, R_i^L$  で表わされる。

$$(\text{領域 I}) = [(R_1^L, R_1^U), (R_2^L, R_2^U), \dots, (R_n^L, R_n^U)]$$

属性として“レコード”（後述）と実セル数を持つ。

- (6) 連結：点  $X$  と領域  $I$  の連結の定義には 4 連結又は 8 連結を用いる。ある点が領域に連結されるとその領域は“再定義”される（図 3-5）。

$$X = (x_1, x_2, \dots, x_n)$$



4 連結の点：○

8 連結の点：○ ⊗

図 3-5 領域と点の連結

1) 4 連結；

$$\begin{aligned} & \exists_i, R_i^L \leq x_i \leq R_i^U \\ & R_j^L - 1 \leq x_j \leq R_j^U + 1 \quad (j \neq i) \end{aligned} \quad (3.9)$$

2) 8 連結；

$$\forall_i, R_i^L - 1 \leq x_i \leq R_i^U + 1 \quad (3.10)$$

(7) 再定義：領域にヒストグラムセル点が連結した時，その領域は再定義される。再定義には連結の仕方によって3通りある。

1) 点が領域 I に外接する時；

(a) その点を含むように領域 I を拡張しても他のクラスタ候補 / クラスタの領域と交差せず，かつ領域 I の“実体積比”が閾値  $R$  以上である場合には，領域 I を拡張する。

(b) それ以外の場合は，その点だけを含む新しい領域 J を作り，領域 I と同じレコードを付ける。

2) 点が領域 I の内側にある時；

領域 I の実セル数に 1 を足す。

(8) 融合：クラスタ，クラスタ候補があるヒストグラムセルを介して連結した場合，連結したものの組み合わせにより以下の3種類の融合が起きる。

1) クラスタ同士の場合；

何もしない（即ち，融合は行なわない）。

2) クラスタとクラスタ候補の場合；

クラスタの核領域の属性を

クラスタ候補の核領域に書き込む。

3) クラスタ候補同士の場合；

処理レベルの新しいクラスタ候補の核領域のポインターを，処理レベルの古いクラスタ候補の核領域を指すように書き換える。

(9) レコード：その領域が属するクラスタ / クラスタ候補の最初の核領域へのポインター及びその核領域が作られた時の処理レベル値（図 3 - 6）。

(10) 処理レベル：

ヒストグラム頻度値，

又は，ヒストグラム頻度値の変化した回数。

(11) 核領域のレコード：クラスタ / クラスタ候補の核領域のレコードのポインターには 0 又はクラスタ番号が書かれている。

(12) 実体積比：領域 I の実体積比  $R$  とは，領域の体積  $V$  と実セル数  $S$  の比で  $R = S/V$  で定める。

(13) 領域選択基準：セルが複数の領域と連結した時に所属させる領域を決める基準

1) その点が内包される領域で，

実体積比が最小の領域。

2) 再定義の際に拡張が可能な領域で，

実体積比が最大の領域。

### 3.2.2 ヒストグラム・モード法のアルゴリズム

ヒストグラム・モード法のアルゴリズムを図 3 - 7，8 に示す。又，1，2 次元の場合の動作例を図

ポインター クラスタ ラベル	処理 レベル値	領域の上限・下限 [ $(R_1^L, R_1^U) \cdots (R_n^L, R_n^U)$ ]
$a_i \rightarrow C_1$	$L_i$	
$a_j \rightarrow O$	$L_j$	
$a_i$	$L_i$	
$a_i$	$L_i$	
$a_j$	$L_j$	

図 3 - 6 領域表の構成

Step-1: 多次元ヒストグラム表を頻度値の降冪の順に並べ換え, スペクトル値の配列  $SV$  と頻度値の配列  $HF$  を作る。  $HF(1)$  は最大の頻度を示す。

Step-2: 処理レベルを  $L_1$  とする。

$$L_1 = HF(1)$$

$SV(1)$  の点を核に領域を作りクラスタ候補とする。領域のレコードは  $[0, L_1]$  となる。

$i = 1$  とする。

Step-3:  $i = i + 1$  として, 処理レベルを更新する。

$$L_i = L_{i+1} - \Delta L$$

Step-4: クラスタ候補のうち処理レベルの値  $1$  が  $1 > L_i + N$

であるものをクラスタとし, その核領域にクラスタ番号をつける。この条件はクラスタ候補を取り巻く谷間の深さが  $N$  以上あることを意味している。

Step-5: 頻度値  $HF(j)$

$$L_i < HF(j) < L_{i-1}$$

のヒストグラムセルに対して Step-6 Step-7 を順次適用する。

Step-6: 現在処理中のヒストグラムセル点  $SV(j)$  とこれまでの処理によって作られたクラスタ候補の領域群との空間的連結関係を調べ;

☆孤立点: どの領域とも連結しない時は点  $SV(j)$  を核に新しい領域を作り, クラスタ候補とする。その処理レベル値は  $L_i$  とする。

☆単連結: 領域  $j$  とのみ連結した場合は, 点  $SV(j)$  を含むように領域  $J$  を再定義する。

☆複連結: 2 つ以上の領域と連結した場合, 領域決定基準に従ってそれらのうち 1 つの領域を決め, その領域を再定義する。連結している領域に異なるクラスタ/クラスタ候補がある場合は, これらのクラスタ/クラスタ候補を融合する。

図 3-7 ヒストグラム・モード法のアルゴリズム

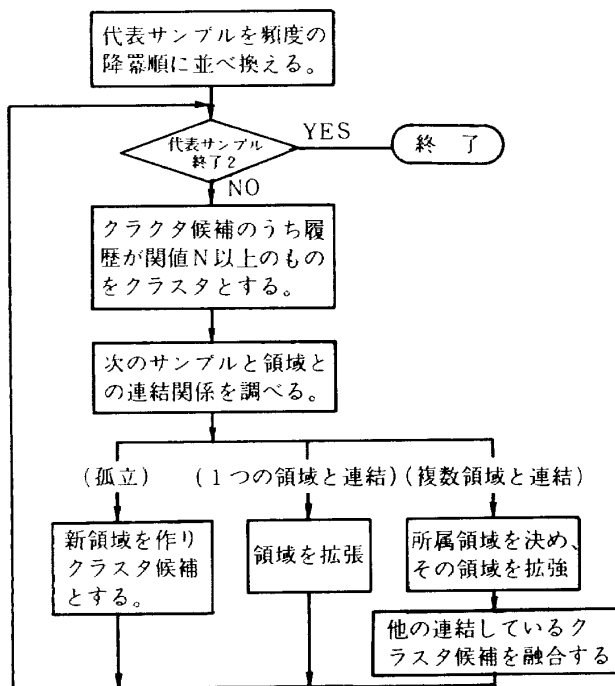


図 3-8 ヒストグラム・モード法のフロー図

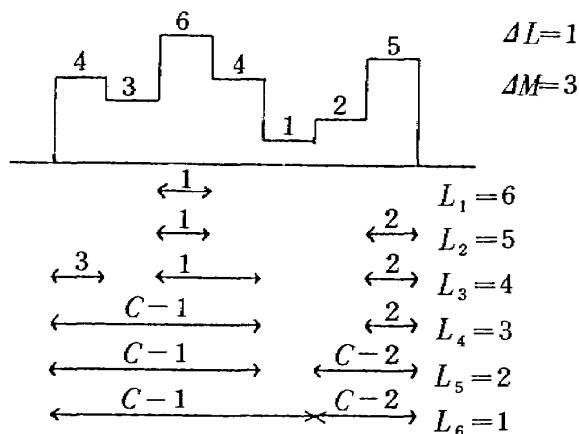
3-9 に示す。

多次元ヒストグラムは 2.2 節で述べた多層ハッシング法で構成されているとする。各ヒストグラムセルは図 2-1(c) のように輝度ベクトルと頻度値が対になっている。

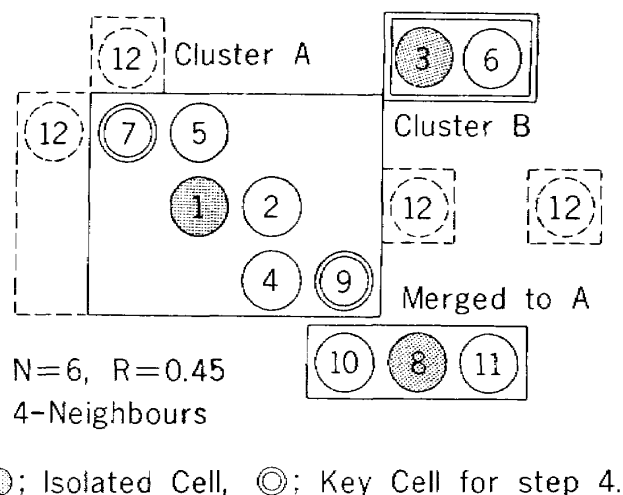
階層的モード法とヒストグラム・モード法を比較すると両手法ともに, 多次元の確率密度の極大点を探索し, 各々の極大点を核にクラスタを順次成長させて確率密度関数のモードによってクラスタ分割をしている点では基本的には同じである。しかし, 処理するデータ量の違いから表 3-1 に示すような違いがある。

解析にかけるデータ量が少ない場合には, ヒストグラムを作成するよりも, 階層的モード法のような形で“データ密度”を定義した方が, 確率密度関数を高精度に表現できる。しかし, 衛星画像のような大規模データの解析処理では,

- (1) ヒストグラムによって高精度に確率密度関数を表わせる,



(a) 1次元ヒストグラムモード法の例



(b) 2次元ヒストグラムモード法の例

図3-9 ヒストグラム・モード法の解析過程例

表3-1 ヒストグラム・モード法と階層的モード法の違い

	確率密度分布関数	クラスタの表現	データ規模	計算負荷 オーダー
階層的モード法	k番目に近い点までの距離	クラスタ内の全ての点のリスト	小規模データ	$N^2$
ヒストグラム・モード法	ヒストグラムの頻度値	多次元直方体の集合	大規模データ	$N'$

(2) 個々のデータ間距離計算荷の点で不可能に近い、

という理由でヒストグラムを用いる方が優れている。

クラスタの表わし方に関しては、クラスタ内の全てのデータのリストを作成する階層的モード法の方が優れている。ヒストグラム・モード法でも同様にする事は可能であるがヒストグラムセル数が $10^4$ 程度であり、現在のアルゴリズムで $10^6$ 程度の距離計算が $10^8$ と約100倍になるので、計算負荷の面で望ましくない。

### 3.2.3 ヒストグラム・モード法のパラメータ

ヒストグラム・モード法で解析結果に影響する主要なパラメータは、次の2つである：

谷間の深さの閾値  $N_v$ ，

実体積比  $R$ 。

(1) 谷間の深さの閾値  $N_v$

極大点は各領域から孤立した点として検出されるが、孤立点が現われるには次の2つの場合がある：

1) 真の極大点、即ち新しいクラスタの核、

2) ヒストグラムの標本誤差による雑音。

問題なのは後者の雑音の場合である。確率密度関数が図3-10(a)のように滑らかであっても、そのヒストグラムには必ず同図(b)のような微小極大点が標本誤差のために現われる。これらの微小極大点も当初は孤立点として検出され、数レベル後に周囲と連結してしまい、微小極大点、即ち標本誤差である事がわかる。従って、全ての孤立点を直ちにクラスタの核とする事は適当でない。

この標本誤差による微小極大点と、分布のモードによる真の極大点を区別するために、ヒストグラム・モード法では“極大点の周りの谷間の深さ”を用いる。ある極大点が標本誤差によるものであれば、その周囲の点との間の頻度の谷間はそれほど深くないであろう。しかし、その極大点が真に分布のモードであれば隣接モードとの間には十分深い頻度の谷間があるはずである。

閾値  $N_v$  は、この谷間の最小の深さを制御する。

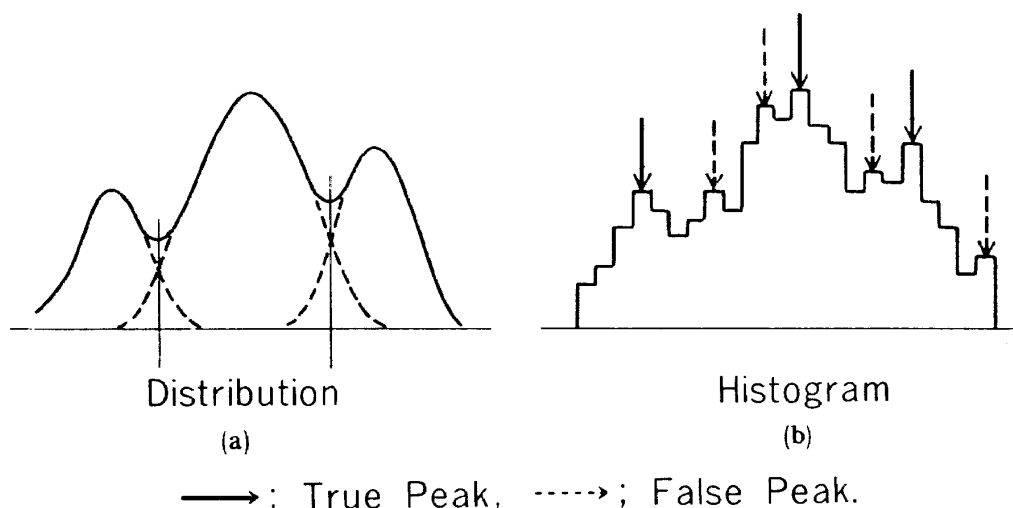


図 3-10 確率密度分布関数とヒストグラム

より大きな  $N_0$  程ヒストグラムの標本誤差に強い解析が可能となるが、あまり  $N_0$  を大きくしすぎると、距離の小さいクラスタ同士が分離できなくなる。

## (2) 実体積比 $R$

ヒストグラム・モード法では計算負荷を軽減するために、クラスタの領域を座標軸に平行な多次元直方体で表現している。このような形での領域の表現は、領域内に含まれた点の近似的表現にしか過ぎないので、領域内の点が座標軸に斜めに分布した場合、実際の点の分布状況と領域の形が掛け離れてくる。

実体積比  $R$  は多次元直方体がどれくらい忠実に実際の状況を表わすかを制御するものである。 $R$  が 1 に近いほど領域と実際の点の分布状況は一致するが、領域の数が大量に必要になり、計算負荷とメモリー量が問題となる。 $R$  が 0 の場合は、最も簡略な解析が行なわれ、最小の計算負荷とメモリーで済む。

## 3.3 クラスタ図の解釈について

衛星画像をクラスタ解析して得られるクラスタ及びクラスタ分類図の持つ意味を考えてみる。

リモートセンシングの分類解析の目的は、1.2 節で述べたように衛星画像を利用目的に応じたいくつかの分類カテゴリーに弁別し、各々のカテゴリーが地上でどのように分布しているか、どれ位存在しているかを把握する事にある。従って、クラスタ解析によって作られるクラスタ分類図が地上のカテゴリーの分布状況を忠実に表現している事が必要である。

しかし、クラスタの処理過程には解析結果をどのような目的に用いるか、という利用上の要求は直接的には反映されない。クラスタ解析では実用的目的を一般化したパラメータ（初期クラスタ、クラスタの数、距離の定義、等）によって間接的に解析の目的を表わすことができるだけである。データ利用のための分類カテゴリーとは全く無関係にクラスタ分類図が作られるとも言えるだろう。

クラスタ解析によって抽出した個々のクラスタが分類上のどのカテゴリーに対応するかは、クラスタ解析だけでは決められない。分類カテゴリーに関するグラントルースデータを用いて、クラスタをカテゴリーに対応付ける解釈処理が、図 3-11 のように、クラスタ解析の後処理として必要である。解釈処理では、

(a) クラスタとカテゴリーの距離関係

(b) クラスタ分類図とトレーニングエリアの画素毎の一致度、

などによってクラスタとカテゴリーが対応付けられる。

この時、クラスタとカテゴリーが 1 対 1 に対応すれば問題ないのであるが、通常は次のような事がしばしば起きる。

(1) 1 つのカテゴリーに複数のクラスタが対応する。

カテゴリー設定が大まかであると、その中に複数のクラスタが抽出される事がある。この場合は殆ど問題ない。

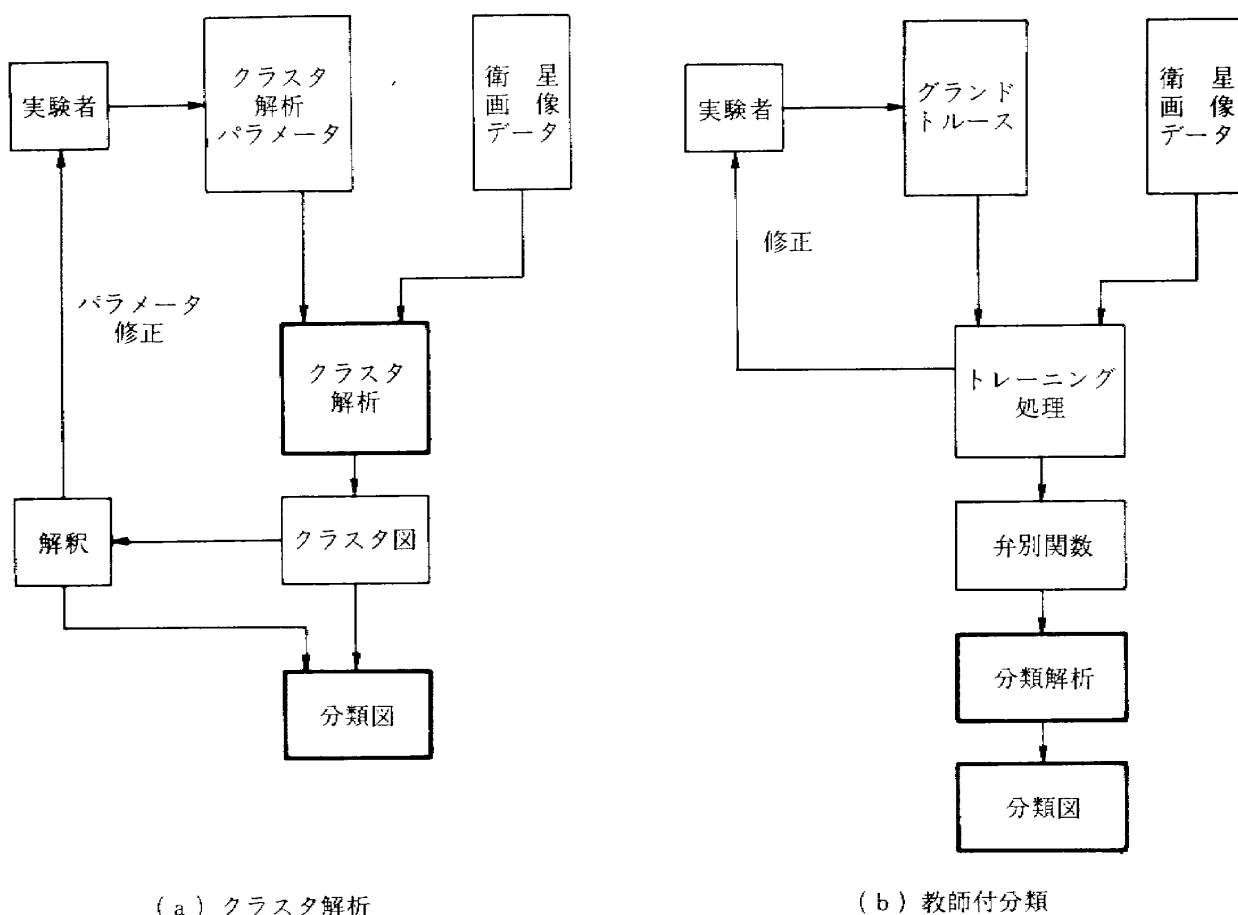


図3-11 分類解析における解釈処理過程の違い

(2) どのカテゴリーとも似ていないクラスタが存在し得る。

分類カテゴリーを十分に設定しない場合に、全てのクラスタをカテゴリーに対応付けようとすると、全く無意味な結果を導く事がある。

(3) 分光的に弁別不加能なカテゴリー設定があり得る。

カテゴリーを設定する場合、それがスペクトルとして弁別可能か否かの検討がなされない事がある。この場合、複数のカテゴリーが1つのクラスタを構成する。

第1の状況は問題ない。一般的にはクラスタ解析は分類しようとするカテゴリーの数よりも多いクラスタを生成するように、パラメータを調節するのが普通である。第2、3の状況は逆に、トレーニングデータの妥当性の検証としてしばしば用いられている。

#### 4. 算法の並列化

ヒストグラム・モード法は、当初から多次元大規模データのクラスタ解析を目的に研究してきたものである。このため、初期においてはFACOM230/75-APUを、後期にはCRAY-1の大型ベクトル計算機を使用し、算法は全面的に並列コードで書かれている。以下、その主要な部分について例示し、衛星画像のような大規模データをベクトル処理する際の考え方について述べる。

ヒストグラム・モード法で大規模な計算負荷が発生する部分は、主に次の3ヶ所である：

- (1) 多次元ヒストグラムを作る際のハッシング計算、
- (2) 多次元ヒストグラム表の頻度順の並べ換え計算、
- (3) ヒストグラムセルとクラスタ領域の連結計算。

以下、各々の場合について述べる。

##### (1) 並列ハッシング

ベクトル計算機では、同時に多数の同じ計算を行なう場合に高い計算効率が発揮できる。このとき、

同時に計算するデータの集まりを計算ベクトルと称する。

衛星画像のハッシング計算で、計算ベクトル長を長くするには、計算ベクトルをスキャンライン方向に設定するとよい。

画素データ  $v$  を式(4.1)で定めると、計算ベクトルは式(4.2)となる。これに対応したハッシュベクトル  $HV$  は式(4.3)で定義される。ランドサット  $MSS$  では  $d=4$  なので1語当たりのbit長に納まるが、 $TM$  のように7バンドもある場合、このままでは計算できない。しかし、式(4.5)~(4.17)に示すアルゴリズムを用いれば、単精度整数演算しかできない計算機でも、式(4.4)でbit長が任意のハッシュ・ベクトルに対してハッシュ・キーを並列に計算できる。

ハッシングの繰り返しループで必要な計算は、式(4.13)とハッシュ法の Step-2 (2.1.2 (2)節)の部分である。これらは付図2に示すようにすれば完全に並列化出来る。但し、付図2において、点線で囲ったベクトル再構成部分は、ベクトル再構成の負荷と処理の完了したデータに対する余分な計算の負荷のバランスによっては適宜取り外す方がよい。

## (2) ヒストグラム表のソーティング

表をあるキーの大きさに順に並べ換えるソーティング計算には、多くのアルゴリズムがあるが、そのいずれもが逐次計算機のための手法で、ベクトル計算機のためのアルゴリズムの報告はない。

そこで、ソーティングのアルゴリズムのうち“入れ替え”によるアルゴリズムをベクトル化したものを付図3に示す。ベクトル化アルゴリズムでは入れ替え位置を探索するステップ及びデータの入れ替えが、完全にベクトル化されるので極めて高速なソーティングが可能となる。

## (3) 連結計算

ヒストグラム・モード法では頻度の順に各ヒストグラムセルと、それまでに作成されたクラスタ(候補)の領域との連結関係を調べ、これによって融合するか、新しいクラスタ候補を作るかを定める。

この計算では、各領域との連結関係の計算負荷が非常に大きい。領域の数は当初1であるが、処理が進むにつれ増加して数百程度にまでなる。これらの

領域と  $10^4$  のヒストグラムセルとの連結関係を調べるので、平均して  $10^6$  回の距離計算が必要になり、並列化によって計算の効率を図らなければならない。

ベクトル計算の長さという点からは、 $v\ell$  個のヒストグラムセルと1つの領域の距離計算の方が効率がよいが、ヒストグラム・モード法のアルゴリズム上この並列化は行なえない。計算は、 $c$  個の領域と1つのセル点の距離計算を同時に行なうという形で並列化する。

付図4に並列コードを示す。(1)のハッシュ計算と同様に、ベクトル再構成の計算負荷と余分な距離計算の計算負荷とのトレードオフが必要である。

$$v_i = \{v_1^i, v_2^i, \dots, v_n^i\} \quad (i=1, \dots, v\ell) \quad (4.1)$$

$$CV = \{v_1, v_2, \dots, v_{v\ell}\} \quad (4.2)$$

$$HV = \{V_1, V_2, \dots, V_{v\ell}\} \quad (4.3)$$

$$V_i = \sum_{k=1}^n v_k^i \cdot \left\{ \prod_{j=k}^{n-1} m_j \right\} \quad (4.4)$$

$$V_i^k = \sum_{\ell=1}^n v_{\ell+4(k-1)}^i \cdot M^{(\ell-1)} \quad (4.5)$$

$$L = \begin{cases} 4 \\ \text{mod}(n, 4) \end{cases} \quad (4.6)$$

$$k = 1, \dots, K$$

$$K = \left\lceil \frac{n+3}{4} \right\rceil$$

$$M = \max_i m_i$$

$$V_i = \sum_{k=1}^M V_i^k \cdot \alpha^{(k-1)} \quad (i=1 \dots v\ell) \quad (4.7)$$

$$\alpha = M^4 \quad (4.8)$$

(ハッシュキーベクトル)

$$HA^j = \{a_1^j, a_2^j, \dots, a_{v\ell}^j\} \quad (4.9)$$

$$HB = \{b_1, b_2, \dots, b_{v\ell}\} \quad (4.10)$$

$$a_i^1 = \text{mod} \left\{ \sum_{k=1}^M \gamma_{k,i}^1 \cdot \beta_k^1, P_1 \right\} + 1 \quad (4.11)$$

表 5-1 クラスタ解析のデータ

	2 次 元 疑似データ	4 次 元 疑似データ	ランドサット 全 体 画 像	ランドサット 補正部分画像
大 き さ	600×100	600×100	3228×2340	1381× 773
次 元 数	2	4	4	4
ク ラ ス タ 数	3	10	?	?
確率密度分布関数	正 規 分 布	正 規 分 布	正規分布?	正規分布 ?
ヒストグラム表	997	9973	19997	9973
データ収集率 (%)	100	97.3	91.4	98.4
解 析 精 度 (%)	99.2	79.6	?	?
Bayes の限界精度	99.3	91.9	?	?

$$b_i = \text{mod} \left\{ \sum_{k=1}^M \gamma_{k,i}^2 \cdot \beta_k^2, P_2 \right\} + 1 \quad (4.12)$$

$$a_i^{j+1} = \text{mod} (a_i^j + b, P_1) + 1 \quad (4.13)$$

$$\beta_1^h = 1, \quad h = 1, 2 \quad (4.14)$$

$$\beta_2^h = \text{mod} (\alpha, P_h) \quad (4.15)$$

$$\beta_k^h = \text{mod} (\beta_{k-1}^h \cdot \beta_2^h, P_h) \quad (4.16)$$

$$\gamma_{k,i}^h = \text{mod} (V_i^k, P_h) \quad (4.17)$$

## 5. クラスタ解析実験・考察

### 5.1 データの記述

ヒストグラム・モード法によるクラスタ解析実験には表 5-1 に示す 4 種のデータを用いた。

#### 5.1.1 疑似データ

アルゴリズムの検証に用いた疑似データは、2 次元及び 4 次元のデータである。いずれも分布関数は正規分布を仮定し、ランダムに与えた平均、共分散を元に Box Muller の正規乱数からデータを生成した (文献 8)。

ランドサットデータをシミュレーションするために各クラスタの分布は適当に重なり合っている (図 5-10, 5-11)。表 5-2 に 2, 4 次元の各クラスタ相互の距離関係を示す。

#### 5.1.2 ランドサットデータ

解析に用いたデータは、Scene ID 1145-00542, 1972 年 12 月 15 日の東海地方のデータである。

図 5-1(a) にシーン全体のバンド 4 の写真を示す。北部の山岳地帯には積雪が見られ、南部の伊豆半島 (右下) 上空には若干雲がかかっているのが見られる。

画像の全体の大きさは 3,240 画素 × 2,340 画素であるが左右に 4 バンド全てのデータが揃わない部分があり、解析に利用できるのは 3,228 画素 × 2,340 画素の領域である。全体の画素数は 7.55

表 5-2 疑似データの分離度 (Separability)

(a) 2 次元疑似データ

	1	2	3
1	—	3.5	3.6
2		—	3.1
3			—

(b) 4 次元疑似データ

	1	2	3	4	5	6	7	8	9	10
1	—	2.8	1.6	2.0	2.7	3.1	2.2	3.3	2.4	3.2
2		—	2.7	2.0	1.8	2.8	3.1	2.4	2.8	2.6
3			—	1.6	2.6	3.7	2.8	3.5	2.6	2.3
4				—	2.3	2.7	2.5	2.9	1.1	2.2
5					—	2.7	2.5	3.7	2.7	3.8
6						—	2.5	2.6	2.0	3.9
7							—	4.0	2.3	4.3
8								—	2.7	2.1
9									—	2.6



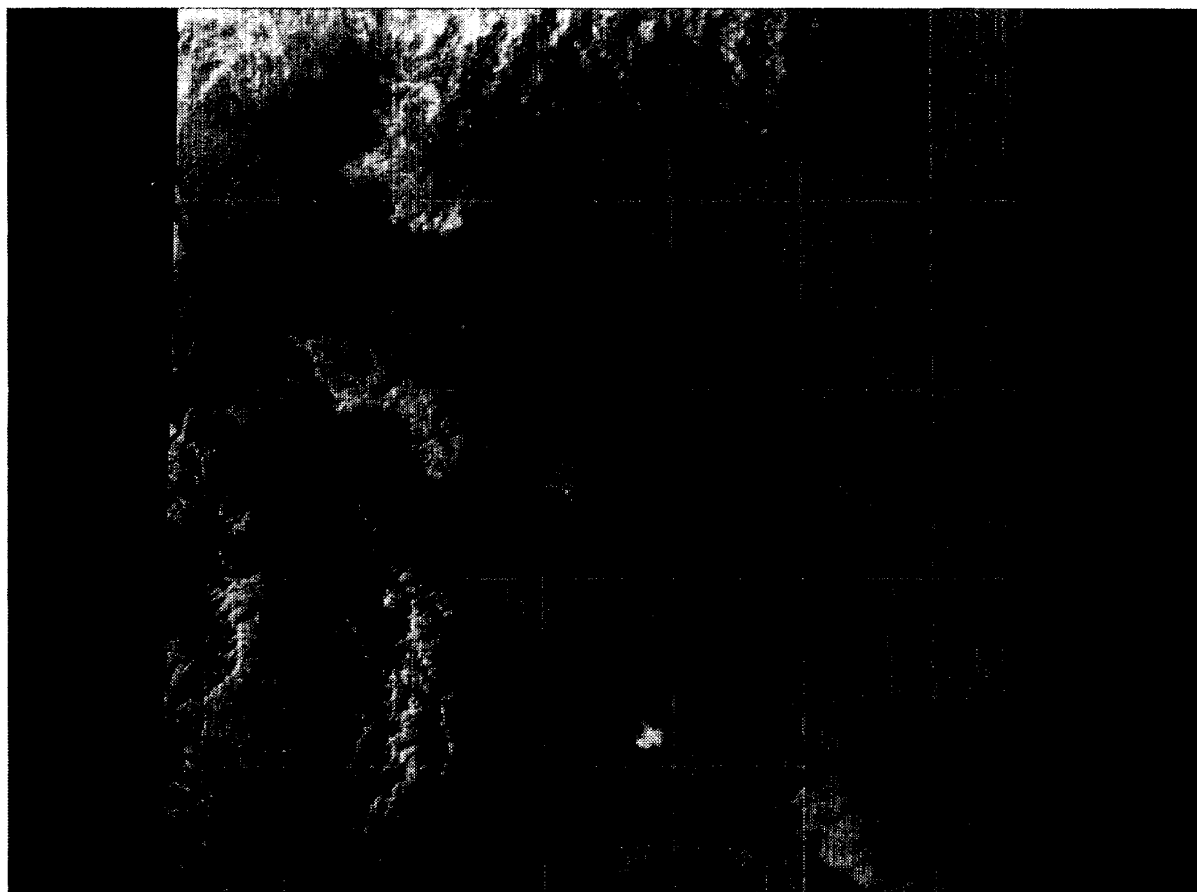


図 5 - 1(a) ランドサット ID 1145-00542 東海地方 1972 年 12 月 15 日 バンド 4

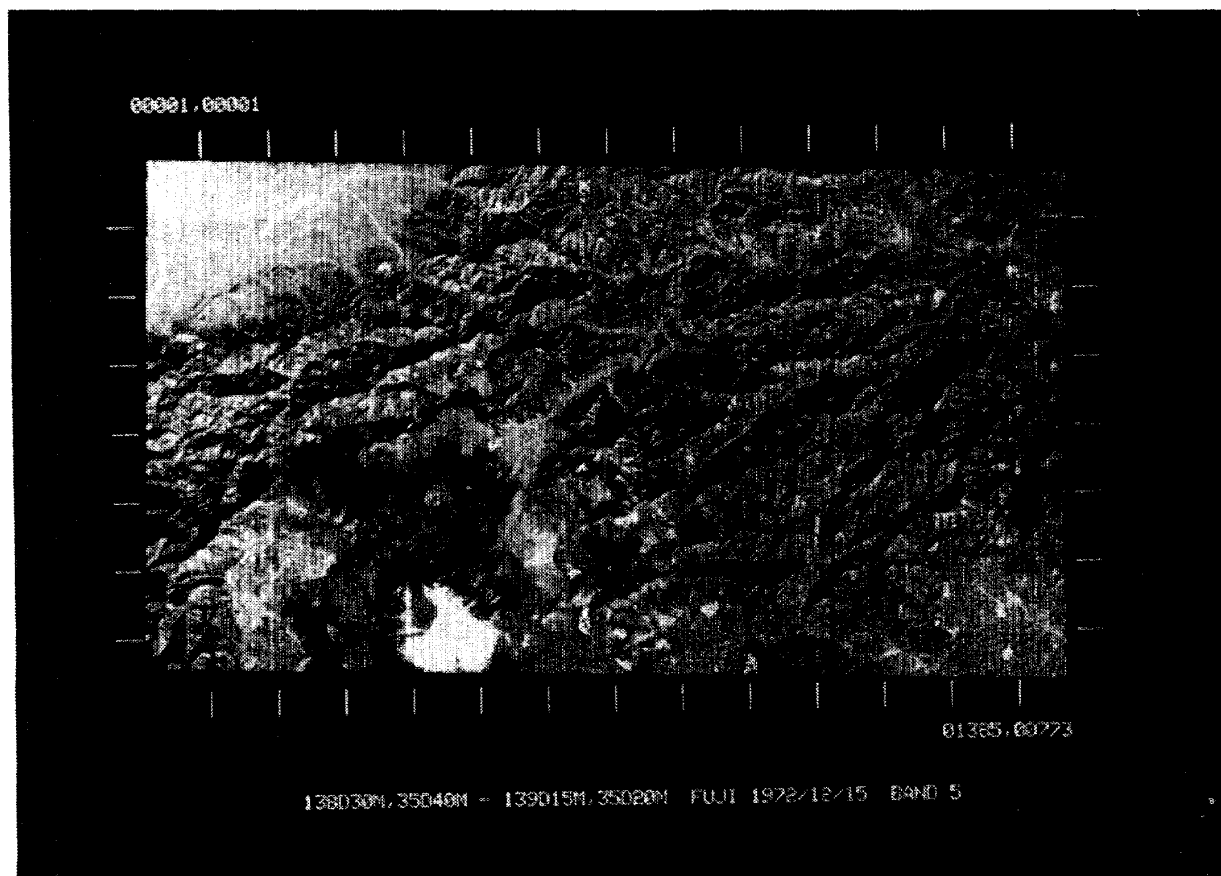


図 5 - 1(b) 地理補正画像 (バンド 5) 東経  $138^{\circ}30' \sim 139^{\circ}15'$  北緯  $35^{\circ}40' \sim 35^{\circ}20'$

$\times 10^6$  になる。

補正部分画像は東海地域全体の画像から東経  $138^{\circ}30' \sim 139^{\circ}15'$ 、北緯  $35^{\circ}40' \sim 35^{\circ}20'$  の地域を抜き出し、幾何学的地理補正を施したもので、 $1,381$  画素  $\times 773$  画素の大きさである。図 5-1 (b) にそのバンド 5 の画像を示す。

## 5.2 多次元ヒストグラムの構成

本節ではランドサットの画像全体のヒストグラム構成に関して述べる。多次元ヒストグラムは、4 章に述べたように並列ハッシング法によって構成している。ランドサット画像の場合、2 スキャンラインを単位に並列化をしており、ベクトル長は  $6,456$  になる。従って、多層ハッシング法でヒストグラム表を更新する操作は 2 スキャンライン毎にしか起動しない。ヒストグラム作成に用いた画像の画素数は  $7.55 \times 10^6$  になる。

### 5.2.1 評価指標の選定及びその閾値

多層ハッシング法ではアルゴリズム Step A (図

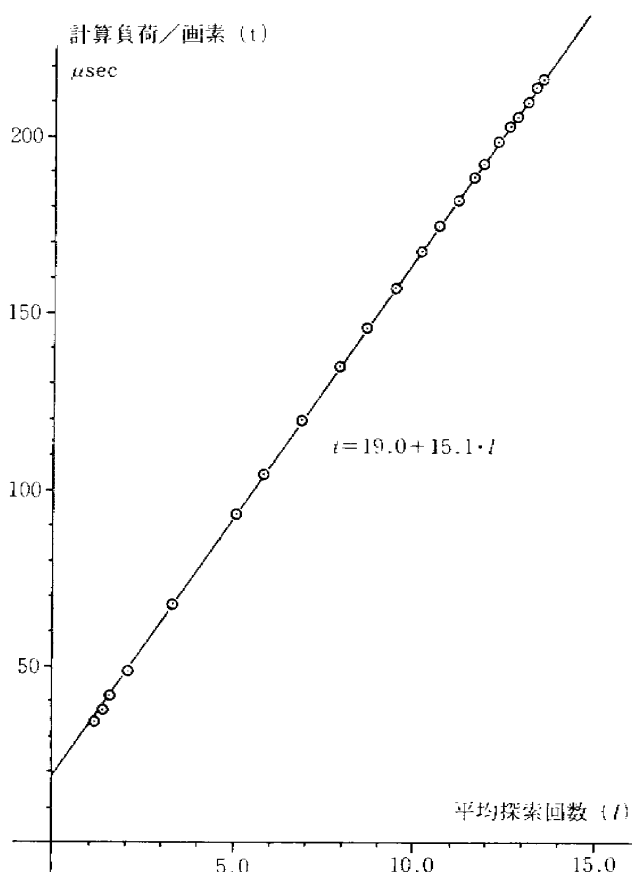


図 5-2 ハッシング法の計算負荷

2-6) で、ハッシングの効率を常に評価しなければならない。ハッシング法の効率は 1 画素当たりの計算負荷で定義できる。この計算負荷は、図 5-2 に示すようにハッシングの探索の平均回数とほぼ完全な比例関係にある。

しかし、Step-A は膨大な原データを取り扱う部分なのでできるだけ計算負荷を減少させることが必要である。平均探索回数は、原データの全ての探索回数が必要なので、計算負荷が大きく、評価指標として適切ではない。このため、より負荷の小さな評価指標として、

- 指標(1) ヒストグラムセルの使用率 (Puse)
- 指標(2) そのヒストグラム表において未登録となった画素の割合 (Prej)

を考える。実際のランドサットデータにおいてこれらの評価指標と平均探索回数との関係を調べたものを図 5-3 に示す。図 5-3 に見られるように指標(2) Prej は平均探索回数と線形関係にある。指標(1) Puse は必ずしも平均探索回数と線形関係にはなく、又、計算効率悪化の始まり近くで急激な変化を示している。従って、指標(2) Prej の方が効率判定には優れている。

図 5-4 からは、ヒストグラム表を更新するための Prej の閾値としては、1% 程度が良いと言える。Prej = 1% はヒストグラム表の 95-98% が使用された状況に相当する。これ以上同じヒストグラム表に登録を続けると、効率が急激に悪化し、この後は処理するスキャンラインの大半のデータが未登録処理となることが見られる。

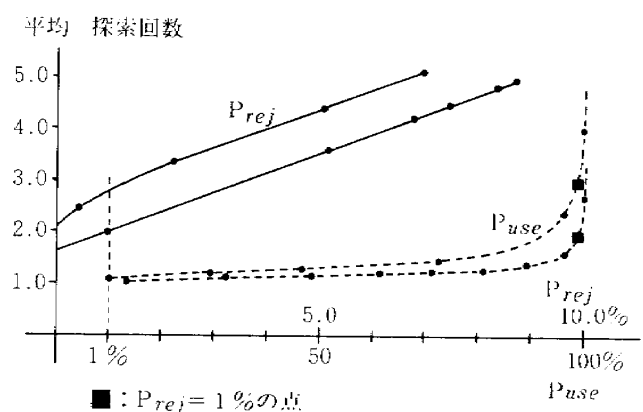
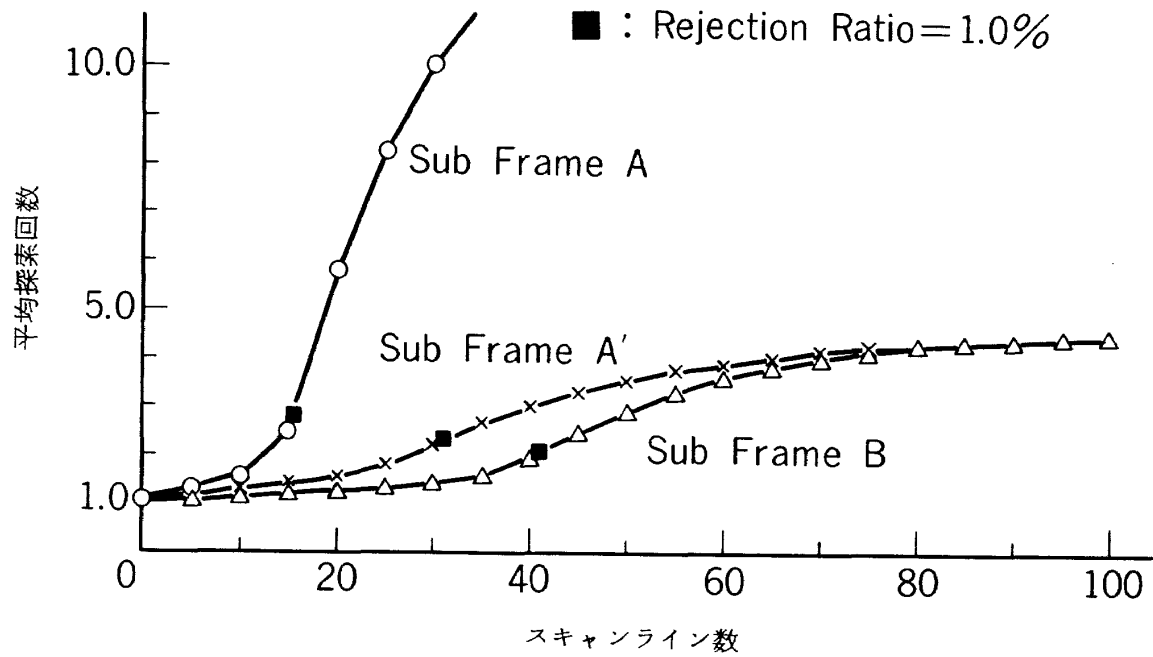
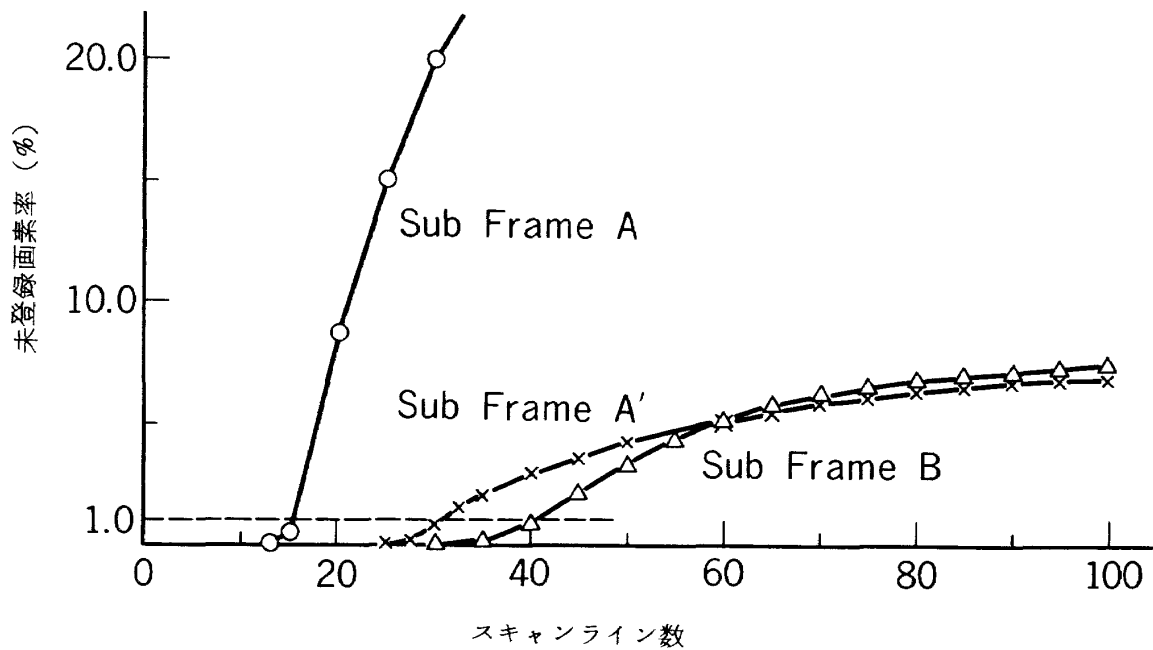


図 5-3 ハッシング法の効率評価指標



(a)



(b)

図 5-4 ハッシング処理の進行による効率の変化

### 5.2.2 多次元ヒストグラムの構成

多層ハッシング法をランドサットデータに適用し、ID.1145-00542 の東海地方のシーン全体から 多次元ヒストグラムを構成した結果を表 5-3 に示す。Step-A, B, C 各々がほぼ等しい計算負荷である。

Step-A では、45 枚のヒストグラム表に全データの 98.8% が登録された。使用したヒストグラムセルの総数は約 92 万個であり、各セルの平均頻度は 8.1 である。各ヒストグラム表毎に処理されたスキャンライン数を図 5-5 に示す。画像の変化の激しい部分（画面北部に相当する部分）ではヒストグラム表の更新が頻繁に行なわれているが、変化の比較的少ない部分ではより多くのスキャンラインが 1 枚のヒストグラム表に収容されている。

92 万個のヒストグラムセルのうちには大量に同じ輝度値を持ったセルが存在しているので、これらを Step-B において融合して 1 つの多次元ヒストグラムとする。Step-B では、それぞれが 98~99% のセル使用率のヒストグラム表を融合するので、衝突が起きる確率は極めて大きい。平均探索回数は図 5-6 に示すように、処理が進むにつれ増加し、最終的には 12.2 回にもなる。Step-A の平均探索回数が 2.0 回であったのに比べると大幅な計算負荷の増大であるが、Step-B で処理する対象は 1/8 にデ

表 5-3 LANDSAT ID.1145-00542 のシーンにおける多次元ヒストグラム作成結果

	基礎的ハッシング法		多層ハッシング法	
	収集率(%)	所要時間(秒)	収集率(%)	所要時間(秒)
Step-A	65.8	1784.7	98.8	331.7
Step-B	—	—	89.6	439.8
Step-C	—	—	91.1	288.5
合計	65.8	1784.7	91.1	1060.0

ヒストグラム表サイズ：19997

ヒストグラム表枚数：45

効率評価開数：Prej=1%

ーク圧縮されたヒストグラムセルなので Step-A とほぼ同等な計算時間で済んでいる。

Step-C においては、Step-A, B でヒストグラム表に未登録となった画素及び、ヒストグラムセルの回収を行なう。理論的には、図 2-7 File B に含まれている画素やヒストグラムセルと同じ観測ベクトルを持つセルが、最終的なヒストグラム表に登録されている可能性があるので、Step-C がないと最終的なヒストグラムが正しい部分ヒストグラムにならない。しかし、表 5-3 に示されるように、そのようなデータの割合は全データの 1.5% と極めて少ない。又、Step-C で処理するデータの大半

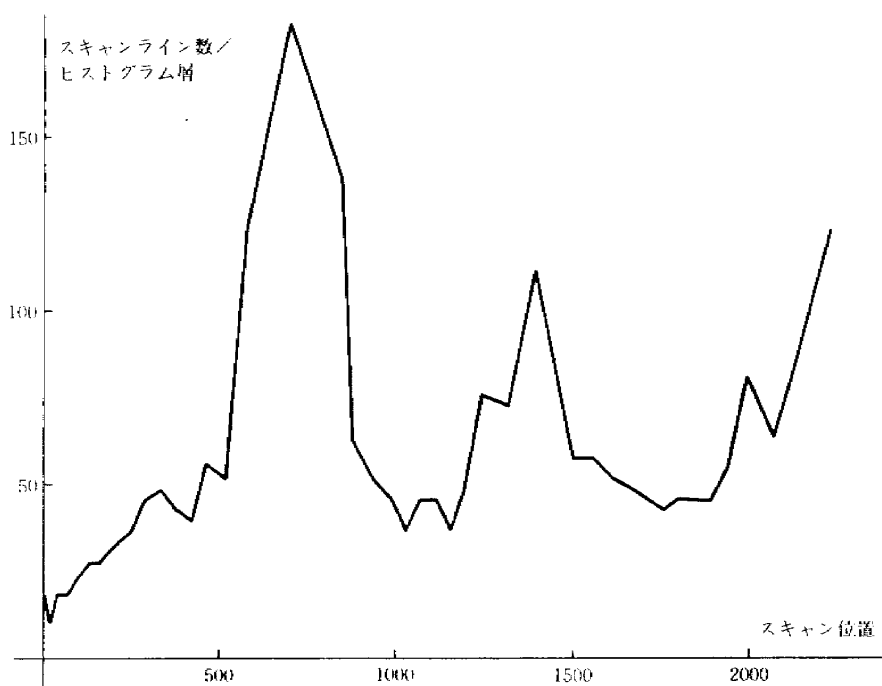


図 5-5 各ヒストグラム表によって処理されたスキャンライン数の変化

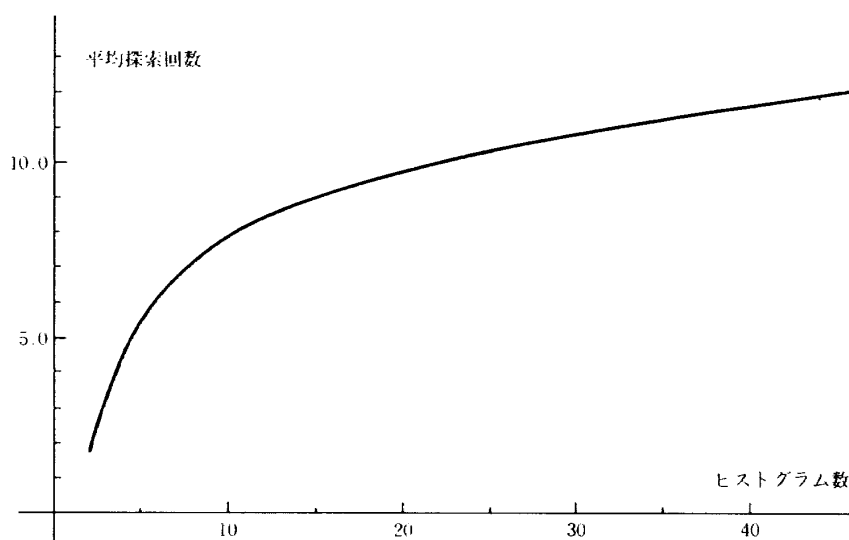


図 5-6 ヒストグラムの融合にともなうハッシュ法の平均探索回数の変動

(8.9/10.4)は結局は未登録になる。この場合、ハッシュングの探索は制限回数まで行なわれてしまうので、計算負荷が収集されるデータの量に比べて相当大きい。従って、コストパフォーマンスの点から Step-C を省き、処理速度を向上させてもあまり問題はない。

### 5.2.3 多層ハッシュング法の効果

本節では多層ハッシュング法の効果を基礎的ハッシュング法と比較して述べる。

表 5-3 には、基礎的なハッシュング法で同じランドサット画像から多次元ヒストグラムを構成した結果を示す。但し、基礎的ハッシュング法ではヒストグラムセルが 100% 使用された際の“再ハッシュング”処理は膨大な計算時間がかかる事が判明したので行っていない。従って、ヒストグラムセル使用率が 100% になった後は新しいセルの登録は全く行なわれない。

表 5-3 に見られるように多層ハッシュング法は基礎的ハッシュング法に比べて、全く同じ大きさのヒストグラム表を用いて、処理時間が 1,785 秒から 1,060 秒に 43% 短縮され、かつデータの収集率は 65.8% から 91.1% へと 25% 以上向上した。

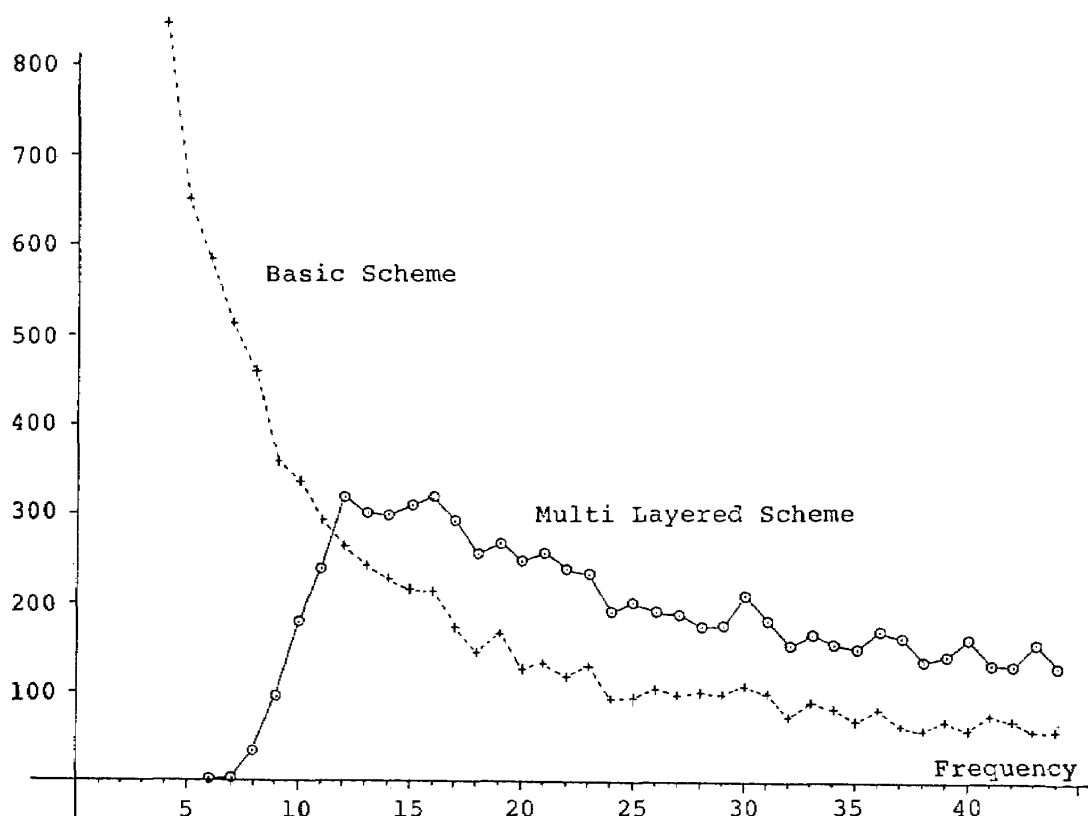
図 5-7 には両手法によって構成された多次元ヒストグラムのセルのうち、最も頻度数が少ない部分の違いを示した。基礎的ハッシュング法では頻度 5 以下のセルが総セルの 35% を占めているが、多層ハ

ッシュング法ではこの部分は少頻度セルとして切り捨てられており、この効率の悪いセルの代わりに頻度 5 以上のセルを余分に集め、ヒストグラムの裾野部分を忠実に表現している。更に、両手法によるヒストグラムの違いは、このような裾野部分のみならず、表 5-4 に示すように、高い頻度の所にも現われており、基礎的ハッシュング法では頻度順位 7 番、9 番という高頻度のセルが見逃されている。

図 5-8 には構成された多次元ヒストグラムに登録されていない画素の割合の、スキャンライン方向の変化を示す。両手法とも、画像の後半を除いて、ほぼ同じような傾向を示し定常的に 20-28% の差がある。しかし、スキャン 1800-2300 にかけて両者の差は著しく拡大している。これは、図 5-1 に見られるように、この近辺から海が画像に現われて来るためである。多層ハッシュング法では各部分画像毎にヒストグラム表を作っているの、海がここで現われた事が Step-A のヒストグラム表に十分反映されるが、基礎的ハッシュング法では既にヒストグラム表が 100% 使用されていて、海の輝度ベクトルを登録できず、全て未登録としてしまっている。更に、海は輝度ベクトルの変化が少ないので、多層ハッシュング法では基礎的ハッシュング法とは逆にこの部分のデータ収集率が他より若干良くなっている。

両手法の大きな差に、使用したランドサット画像に依存している所もある。図 5-9(a)に見られるように、データの最初の 100 スキャン程は非常に輝

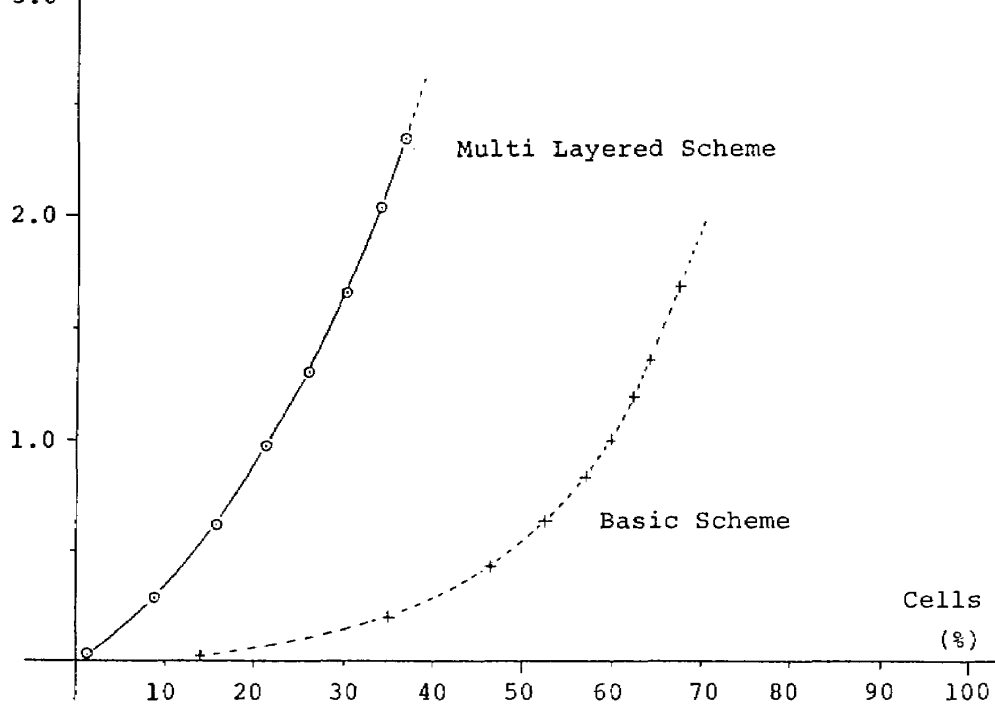
Number of Cells



(a) 頻度の小さい部分での同頻度セル数

Gathering

(%)



(b) 微小頻度のセルのヒストグラム表占有率

図5-7 多次元ヒストグラムの裾野部分における多層ハッシング法と基礎的ハッシング法の違い

度の変化が激しくヒストグラムが作り難い画像であった。そのため、初期段階で出来上がったヒストグラム表の低効率がそれ以後に波及し、全体としての効率を低下させたとも言える。

図5-9に、部分画像のヒストグラムを基礎的ハッシング法と多層ハッシング法で作った場合に、未登録となった画素の割合の変化、計算時間、及びヒストグラム作成の効率を示す。部分画像は図5-1(a)の画像から150ラインずつ帯状に切り出して作

表5-4 LANDSAT ID.1145-00542のシーンにおける最頻輝度ベクトル例

順位	頻度値	輝度ベクトル
1	17267	(12, 6, 3, 1)
2	16497	(21, 17, 18, 8)
3	14769	(21, 17, 19, 9)
4	12721	(12, 4, 3, 0)
5	12717	(12, 6, 4, 1)
6	12125	(21, 17, 18, 9)
7	12087	(17, 8, 4, 0) **
8	12071	(14, 6, 4, 1)
9	12053	(17, 8, 4, 1) **
10	10983	(21, 17, 20, 10)

\*\*：基礎的ハッシング法では検出されなかった輝度ベクトル

った。未登録画素率と所要計算時間の積でヒストグラム構成の効率を評価すると、同図(c)に見られるように殆どの場合、多層ハッシング法の方が高効率でヒストグラムを作っていると言える。

以上の事から多層ハッシング法について以下の事が言える。

- (1) 画像に変化の激しい部分が存在する場合、多層ハッシング法の効果が顕著になる。
- (2) 多層ハッシング法では、クラスタが画面のどこにあっても十分な割合で存在すれば、必ず多次元ヒストグラムに反映される。
- (3) 基礎的ハッシング法ではヒストグラムの裾野の部分が十分表現されていないのみならず、頻度の大きなセルにおいても見落としがある。多層ハッシング法にはこの欠点はない。
- (4) 画像の変化が少ない所では両手法とも大きな差はないが、多層ハッシング法の方が計算時間の長い分だけ結果も良くなる。

### 5.3 解析実験結果及び考察

ヒストグラム・モード法によるクラスタ解析実験を5.1に述べた4種類のデータに対して行なった。

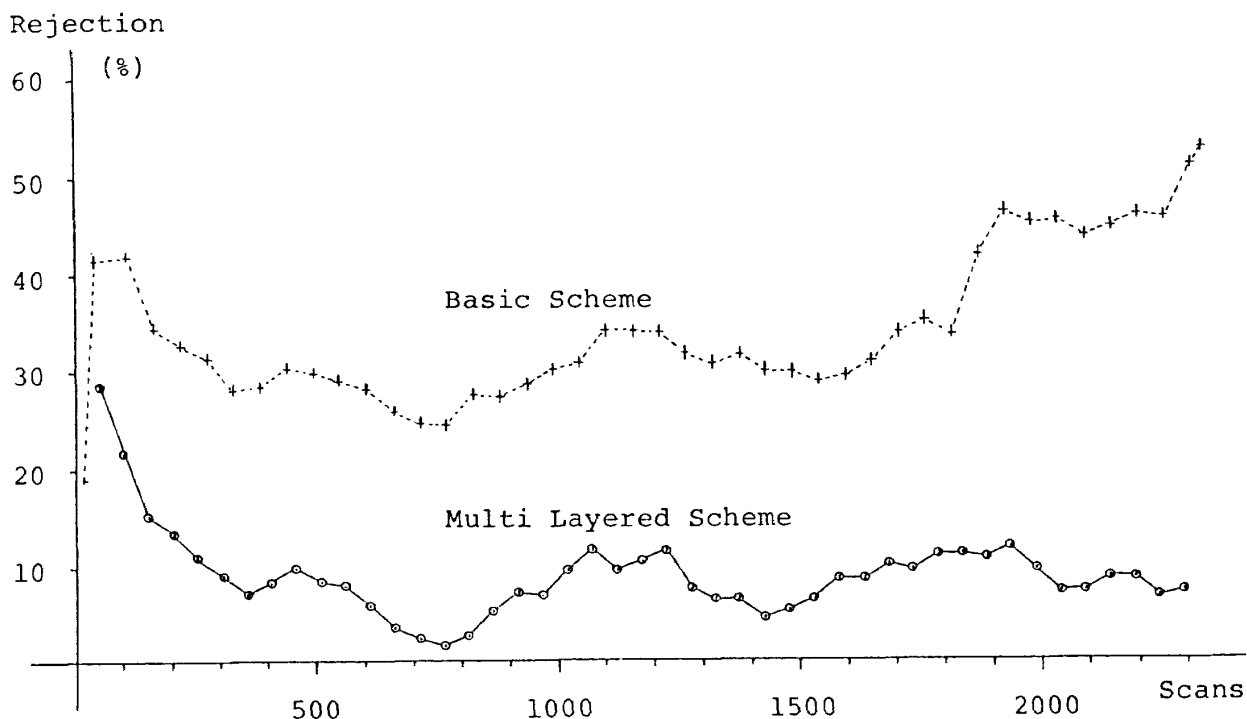


図5-8 ヒストグラムに未登録な画素の各画像位置における違い

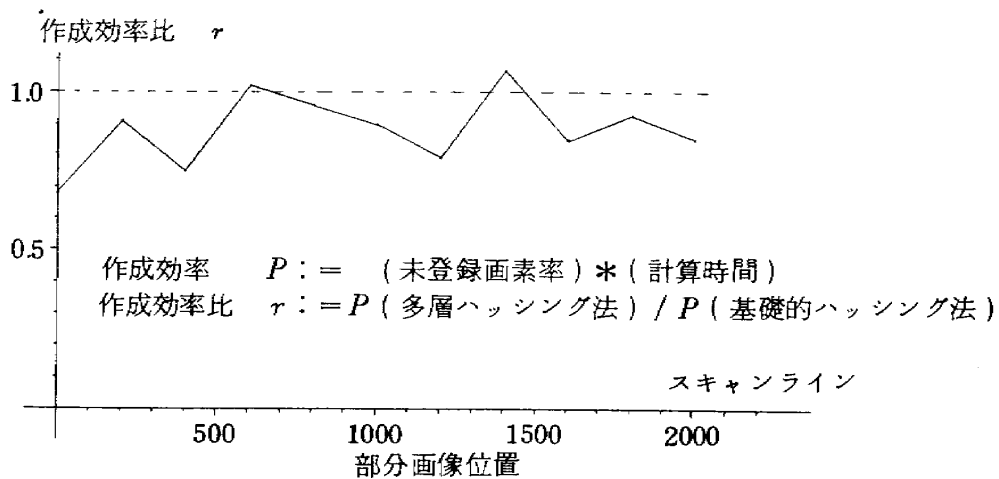
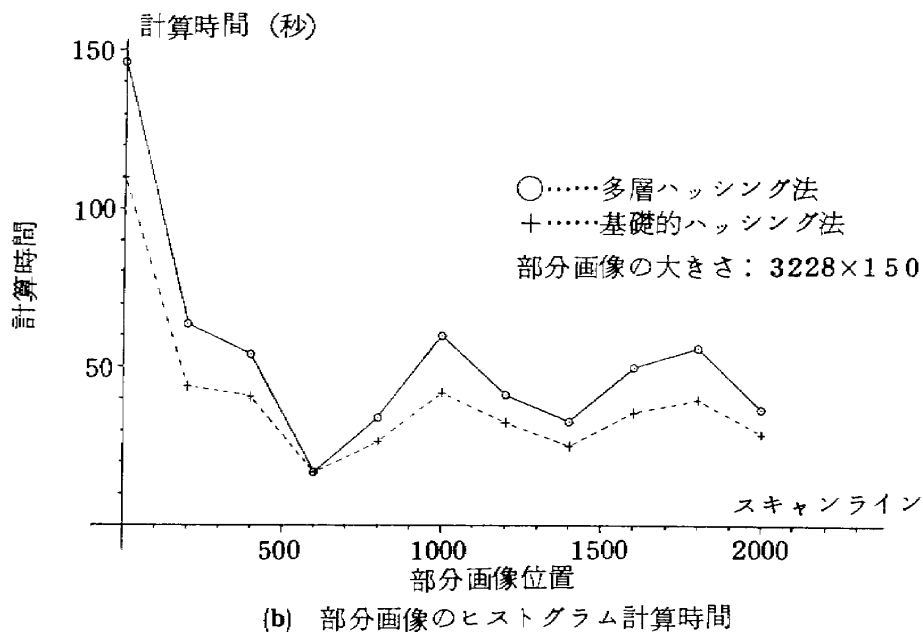
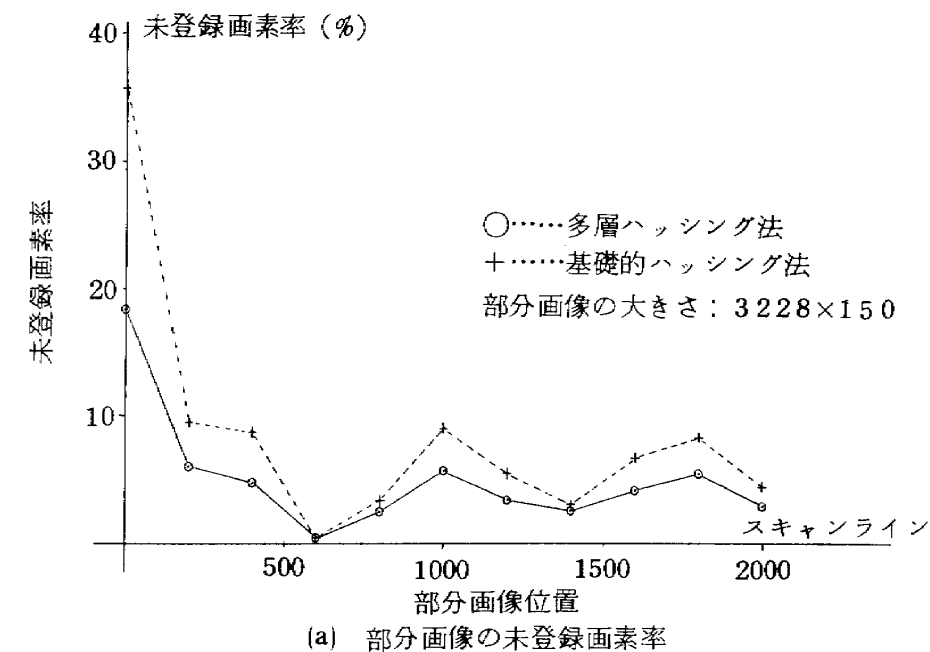
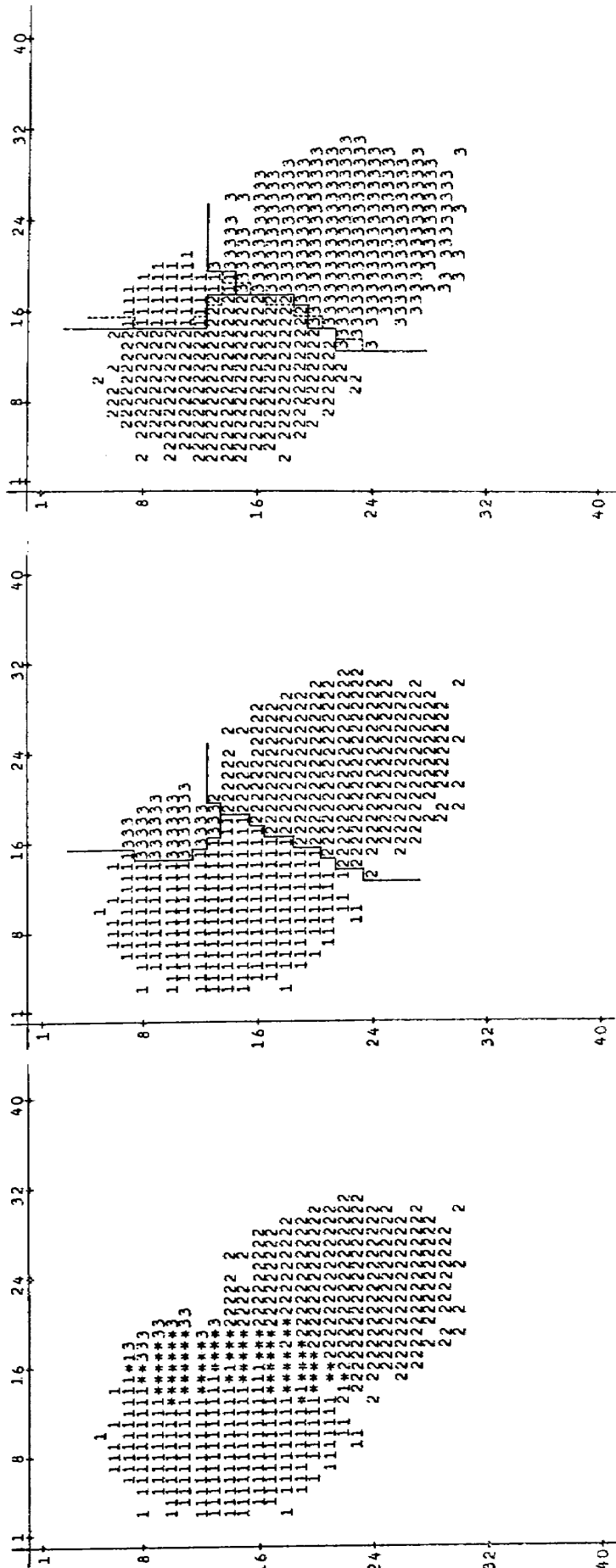


図5-9 多層ハッシング法の効果



表5-5 LANDSAT ID1145-00542.  
トレーニング・クラスの誤差行列

	ト レ ー ニ ン グ ク ラ ス	1	2	3	4	5	6	7	8	9	10	11	12	13	14	15	16	17	18	19	20	21	22	23	24	25	26	27	28	29	30	31	32	33	34	35	
1市街地	226	863	19												32	33																					
2日陰	455	879																																			
3耕地	262	35	656	22	15																																
4耕地	45		19	778	08																																
5耕地	133		15		805	16																															
6果樹園	63				15	921																															
7湖沼	1	72				889																															
8湖沼	2	1486				863																															
9雪	1	153																																			
10雪	2	77																																			
11雪	3	30																																			
12雪	4	72																																			
13雪	5	20																																			
14河川敷	31	04																																			
15河川敷	30																																				
16草地	1	50																																			
17草地	2	145																																			
18草地	3	35																																			
19草地	4	40																																			
20岩石地	150																																				
21ササ地	25																																				
22新植地	18																																				
23裸地	205																																				
24カラヤ	114																																				
25カラヤ	108																																				
26カラヤ	60																																				
27カラヤ	227																																				
28カラヤ	44																																				
29カラヤ	115																																				
30カラヤ	55																																				
31カラヤ	40																																				
32カラヤ	30																																				
33カラヤ	42																																				
34カラヤ																																					
35カラヤ																																					
36カラヤ																																					



(c) ヒストグラム・モード法によるクラス  
タ境界 (点線は Bayes 境界)

(b) Bayes の最適弁別境界

(a) 2 次元疑似データ分布図

図 5-10 2 次元疑似データのクラスタ解析

### 5.3.1 解析結果の解釈

クラスタ解析結果は3.3節で述べたように何らかのランドトールスデータと突き合わせ、解釈をしなければ意味ある結果にならない。

以下に述べる解析実験では、この解釈は全てトレーニングクラスとの距離 (Separability式5.1) によって決定した。

$$d = \frac{B_{12}^t (\mu_1 - \mu_2)}{(B_{12}^t K_1 B_{12})^{1/2} + (B_{12}^t K_2 B_{12})^{1/2}} \quad (5.1)$$

$$\text{但し } B_{12} = \left( \frac{K_1 + K_2}{2} \right)^{-1} \cdot (\mu_1 - \mu_2)$$

クラスタにラベルを付けるため、各クラスタ分類されたデータから平均、共分散を求め、トレーニングクラスとの距離を計算し、最短距離のクラスのラベルをそのクラスタのラベルとした。

疑似データの解析結果の解釈には、合成した疑似データから計算したトレーニングクラスの平均・共分散を用いた。これは真の平均・共分散の一致推定量となっている。

ランドサットデータに対しては、真のクラス・平均・共分散全てが未知であるから、教師付き分類と同様に実験者によりトレーニングエリアを選定し、トレーニングクラスを作成した。表5-5に選定したトレーニングクラスの一覧及びその誤差行列 (Confusion Matrix) を示す。

### 5.3.2 疑似データの解析結果

#### (1) 2次元疑似データ

図5-10(a)には、実験に用いた2次元データの分布図を示す。図中、数字1, 2, 3はそれぞれのクラスタ分布を示し、\*印の所は分布が重なっている部分を示す。同図(b)にはこのデータの確率密度関数が全て既知であるとした場合、Bayesの最適弁別境界を示す。分布に重なりがあるため、いかなる弁別境界を設定してもこれ以上の高精度は達成できない。

このデータをヒストグラム・モード法で解析した。谷間閾値  $N=2$ 、実体積比  $R=0.3$  で解析した。15の領域が使用され、3個のクラスタが抽出された。

表5-6 クラスタ vs 真のクラスの Separability 距離 (2次元疑似データ)

クラスタ	真のクラス		
	1	2	3
1	0.009	2.251	> 3.0
2	2.327	0.011	2.914
3	> 3.0	2.937	0.004

表5-7 クラスタ vs 真のクラスの 誤差行列 (2次元疑似データ) (%)

クラスタ	真のクラス		
	1	2	3
1	99.3	0.2	1.7
2	0.3	99.8	0.0
3	0.4	0.0	98.3

弁別精度：99.15%

解析の結果図5-10(c)に示す境界が設定された。図中、点線はBayesの最適弁別境界である。その違いは12点程度であり、かつ境界のズレは±1に止っている。表5-6には抽出されたクラスタと真のクラスの一致度をクラスタ間距離で示す。又、表5-7には解析結果の誤差行列を示す。

2次元データの解析に関してはヒストグラム・モード法は十分な解析性能を持っていると言えるであろう。

#### (2) 4次元疑似データ

ランドサットMSSデータを模擬した4次元データの解析を試みた。図5-11(a)には実験に用いた4次元データの各クラスタの1σ境界面の2次元平面への投影図を示す。

全データ60,000の97.4%が9,973個のヒストグラムセルに集約された。この4次元ヒストグラムを谷間閾値  $N=2$ 、実体積比  $R=0.3$  で解析した所139の領域を使用して、15のクラスタが抽出された。15個のクラスタのうち5個は要素数1の微小クラスタなので無視した。

得られたクラスタの1σ境界投影図を図5-11(b)に示す。全体の構成はよく似ており、元のクラスの形をよく推定している事が認められる。この結果を

表5-8 クラスタ vs 真のクラスの Separability 距離  
(4次元疑似データ)

クラスタ	真 の ク ラ ス									
	1	2	3	4	5	6	7	8	9	10
1	<u>0.16</u>		1.44	1.76			2.20		2.45	
2		<u>0.17</u>		1.96	1.57					
3	1.71		<u>0.10</u>	1.40					2.36	2.06
4		2.37	2.03	<u>0.51</u>					0.67	2.21
5	2.29	1.96	2.35	2.40	<u>0.17</u>		2.29			
6				2.45		<u>0.08</u>	2.30		1.86	
7	2.20						<u>0.06</u>		2.25	
8		2.18						<u>0.05</u>		2.12
9	2.36	2.14	2.01	0.71	2.46	2.16	2.42		<u>0.49</u>	
10			2.43	2.31				2.06		<u>0.06</u>

但し、空白は3.0以上を示す。

表5-9 クラスタ vs 真のクラスの誤差行列  
(4次元疑似データ) (%)

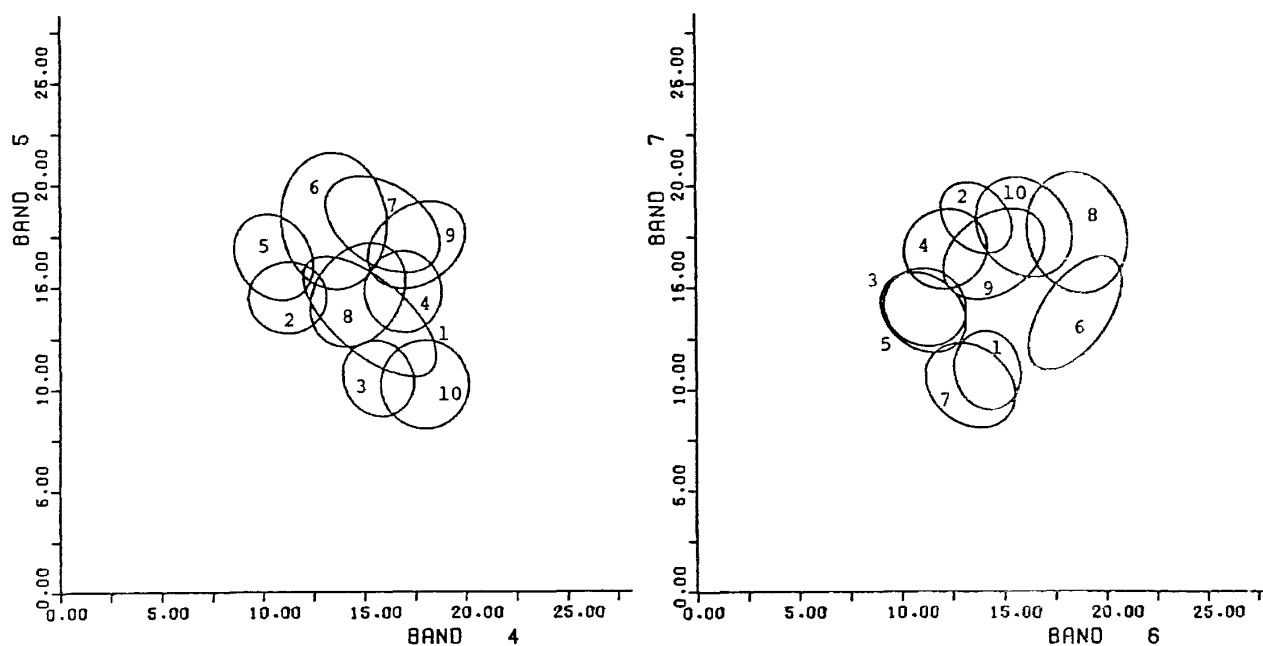
クラスタ	真 の ク ラ ス									
	1	2	3	4	5	6	7	8	9	10
1	<u>81.0</u>	0.0	2.8	0.0	3.4	0.1	2.6	0.0	0.7	0.0
2	0.1	<u>85.0</u>	0.4	0.7	2.1	0.2	0.0	4.4	2.1	1.0
3	11.5	0.1	<u>80.1</u>	2.3	2.1	0.0	0.1	0.0	1.1	0.3
4	4.5	3.6	11.9	<u>49.2</u>	0.3	0.8	0.3	0.1	33.7	0.6
5	0.2	10.5	0.3	0.0	<u>89.8</u>	0.4	0.3	0.0	0.2	0.0
6	0.0	0.3	0.0	0.0	0.5	<u>90.0</u>	0.3	0.6	1.9	0.0
7	2.0	0.0	0.1	0.0	1.9	2.2	<u>94.6</u>	0.0	1.0	0.0
8	0.1	0.2	0.1	0.2	0.0	0.8	0.0	<u>91.7</u>	0.4	3.9
9	0.5	0.2	0.4	45.7	0.0	5.4	1.8	0.9	<u>58.4</u>	0.1
10	0.1	0.1	3.9	1.8	0.0	0.0	0.0	2.4	0.6	<u>94.1</u>

弁別精度：79.59%

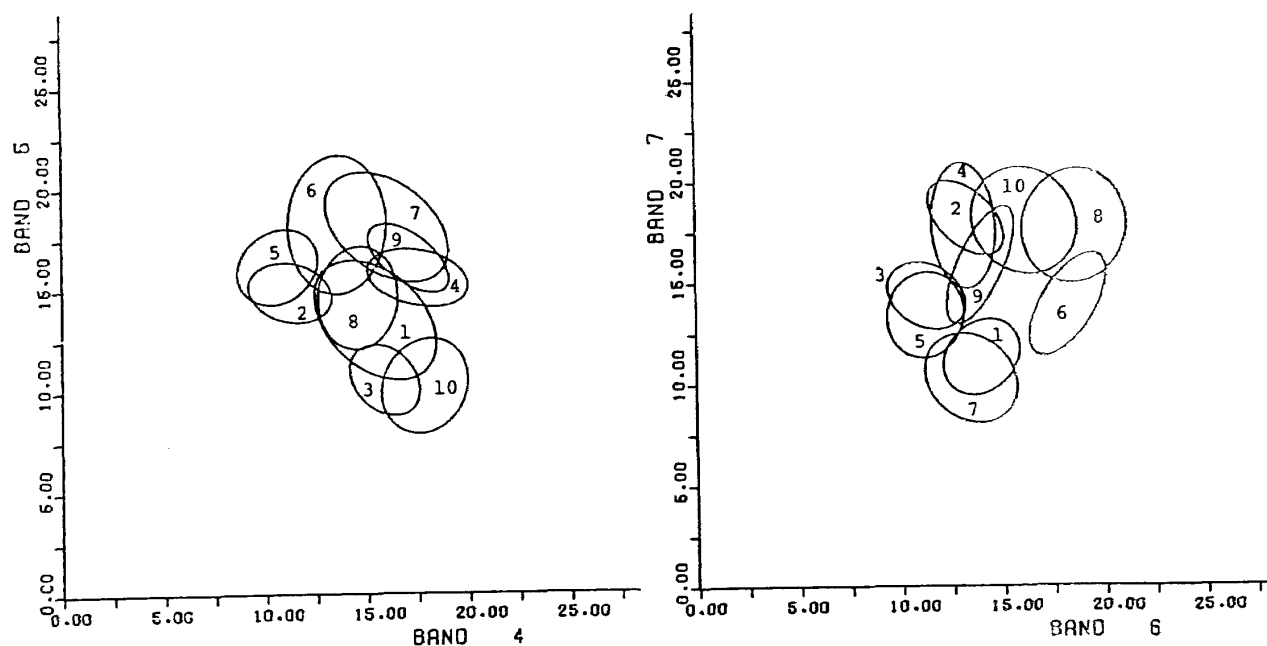
表5-10 最尤法解析結果 vs 真のクラスの誤差行列  
(4次元疑似データ) (%)

最 尤 法 解析結果	真 の ク ラ ス									
	1	2	3	4	5	6	7	8	9	10
1	92.5	0.0	3.6	1.3	0.1	0.0	1.0	0.0	0.4	0.0
2	0.0	95.4	0.0	1.9	3.7	0.2	0.0	0.6	0.0	0.4
3	4.2	0.0	90.0	5.2	0.3	0.0	0.1	0.0	0.0	0.8
4	1.3	1.4	5.1	77.1	0.7	0.1	0.2	0.1	13.7	1.0
5	0.3	2.2	0.2	0.8	94.3	0.2	0.6	0.0	0.1	0.0
6	0.0	0.1	0.0	0.0	0.1	97.2	0.2	0.3	1.6	0.0
7	1.3	0.0	0.0	0.3	0.6	0.6	96.8	0.0	0.7	0.0
8	0.0	0.6	0.0	0.0	0.0	0.4	0.0	97.2	0.3	1.6
9	0.3	0.1	0.0	12.3	0.1	1.4	1.0	0.2	82.9	0.1
10	0.0	0.2	1.1	1.3	0.0	0.0	0.0	1.6	0.3	96.0

弁別精度：91.91%



(a) 4次元疑似データ各クラスターの境界投影図



(b) ヒストグラム・モード法によって得られた各クラスターの境界投影図

図 5-11 4次元疑似データのクラスター解析

クラスター間距離で評価した結果を表 5-8 に示す。又、解析結果の誤差行列及び Bayes の最適弁別での誤差行列を表 5-9, 5-10 に示す。

他のクラスとよく分離している 5, 8, 10 のクラスターに関してはほとんど限界に近い解析精度が得られたが、他のクラスと分離の悪いクラスター 4, 9 などでは解析精度は悪い。しかし、全体としては Bayes

の上限 91.9% に対して 79.6% の分類精度が得られており、クラスターを十分正しく推定できたと言えるであろう。

ヒストグラム・モード法によってクラスター 4, 9 が分離しにくい原因を考えてみる。クラスター 4 と 9 の分布の重なり程度が問題であるが、分離度は Separability で 1.1 (表 5-2 (b)) であり、これ

は両クラスが大分重なっている事を意味する。簡単のために、1次元の場合に分離度 1.1 の2つの分布の重なり具合をプロットしたものを図5-12(a)に示す。同図(b)には分離度 2.1 (即ちクラス10の最小分離度) の場合を示す。ヒストグラム・モード法の基本的手法が、混合確率密度関数の極大点・谷間の探索であることを考えると、同図(a)のような状況はヒストグラム・モード法の限界に近く、クラスタ推定精度が低下した事は当然と言える。又、同図(b)のように混合確率密度関数が非常に明瞭に分離している場合に ( クラスタ 5, 8, 19 ), ヒストグラム・モード法では Bayes の限界に近い推定精度が得られた事は、本手法の正しさを証明していると言えるであろう。

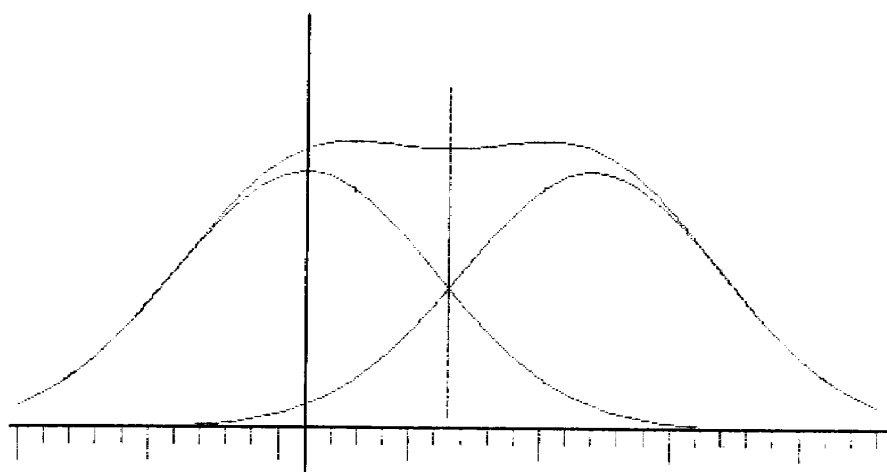
### 5.3.3 ランドサットデータの解析実験

クラスタ解析は東海地域の全体及びその補正部分画像に対して行なった。

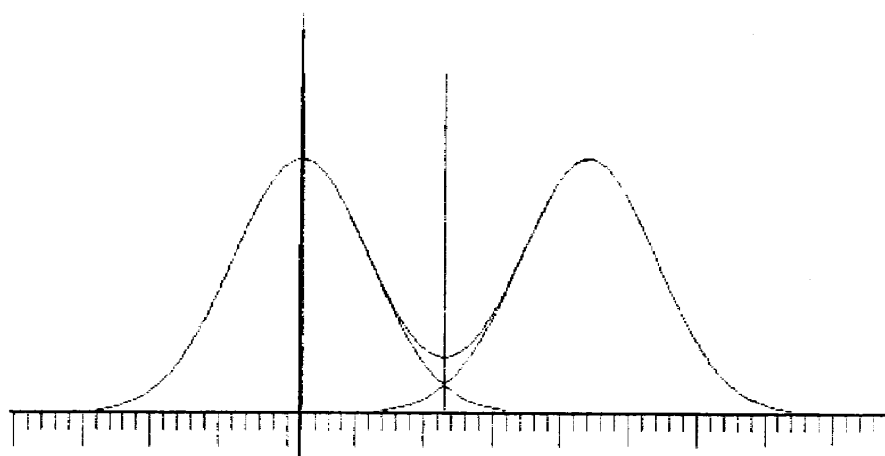
ヒストグラムはシーン全体に対しては 19,997 個のヒストグラムセルを用い、補正部分画像に対しては 9,973 個のセルを用いて構成した。ヒストグラム構成の結果を表5-11に示す。

構成したヒストグラムにヒストグラム・モード法を適用し、前述したトレーニングクラスによってクラスラベルをつけた部分画像の分類結果を図5-13に示す。又、同じトレーニングデータによって最尤法分類した結果を図5-14に参考として示す。両図を目視によって比較すると、次のような違いがみられる。

(1) ヒストグラム・モード法の解析結果の方が、山



(a) Separability = 1.1



(b) Separability = 2.1

図5-12 混合確率密度関数 ( 1次元, 等分散, 等先験確率 )

表 5-11 LANDSAT 1D.1145-00542 のシーン  
における多次元ヒストグラム作成結果

	全 体 画 像 収集率(%) 所要時間(sec)		補 正 部 分 画 像 収集率(%) 所要時間(sec)	
Step - A	98.8	331.7	99.1	32.3
多 層 Step - B	89.6	439.8	98.2	5.1
ハッシング法 Step - C	91.1	288.5	98.4	3.3
合 計	91.1	1060.0	98.4	40.7
ヒストグラム表の長さ	19997		9973	
層数	45		2	
処理時間/画素 ( $\mu\text{sec}$ )	140		56	

岳地帯に谷間が入り込んでいる様子がよく分かる。

- (2) 観測画像全体にわたって、ヒストグラム・モードによる画素の空間的まとめあげが見られ、最尤法分類のような個々の画素単独の分類による解析精度の劣化が抑制されている。

図 5-13 の結果は谷間の閾値  $N_v$  が 1 の時、即ち最も細かく分割した時の結果である。 $N_v$  を大きくすると、互いに近いクラス同士から段々に融合していき解析結果も変動して来る。谷間の閾値  $N_v$  1, 5, 10, 50 と変えて、ヒストグラム・モード法を実行し、画素数の多い主要なクラスについて、 $N_v$  が大きくなるにつれクラスが統合されて行く状況のデンドログラムを表 5-12(a), (b) に示す。

$N_v = 50$  という大きな閾値を設定すると補正部分画像、全体画像ともに、クラス同士の連結 (チェーンニング) が進み、クラスにわたるクラスの統合が起きて来る。連結されるクラスのラベルを調べると、 $N_v = 1, 5, 10$  では異なったクラスのラベルを持っていたクラス同士が、 $N_v = 50$  で 1 つのクラスに融合する事が、

〔草, カラマツ, 耕地, 日陰林〕

〔日陰, 日陰雪, 湖, 溶岩地, (日陰林)〕

などのグループ内で見られる。これは、これらのグ

ループに対して設定したトレーニングクラス間の分離度が低く、それに対応したクラス間に  $N_v = 50$  という大きな谷間がない事を意味している。

$N_v = 5$  程度では多くのクラスは  $N_v = 1$  の時と殆ど同じ構成で同じクラスラベルを持つ。クラスの混合が見られたクラスが 1 つあるが、この場合も  $N_v = 1$  でのクラス決定での最短距離が 1.54 もあり、元々のクラス決定に無理があったものである。

$N_v = 10$  の場合も  $N_v = 5$  の場合によく似ているが、クラスの融合が更に進み、クラスにわたる混合が明らかなクラスが 1 つあった。

以上の事から、ランドサットデータにおける谷間の閾値  $N_v$  について；

- (1)  $N_v$  を大きくすると、互いに近いクラス同士が連結されるチェーン効果が起きる。
- (2) 設定したトレーニングクラスに対しては  $N_v = 50$  は大き過ぎる。特に〔草, カラマツ, 耕地, 日陰林〕及び〔日陰, 日陰雪, 湖, 溶岩地, (日陰林)〕などの間は確率密度の谷間がそれほど大きくない。
- (3)  $N_v = 1 \sim 5$  の間は余り変わらない。  
などがみとめられる。

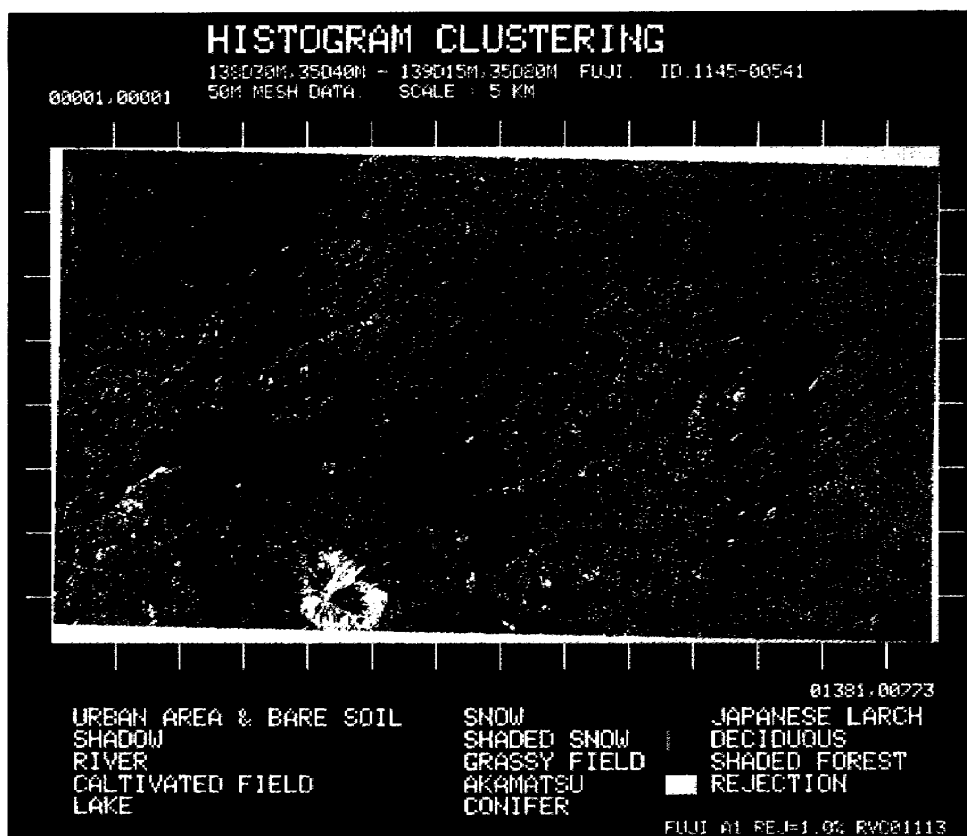


図 5-13 クラスタリング解析による植生分類図

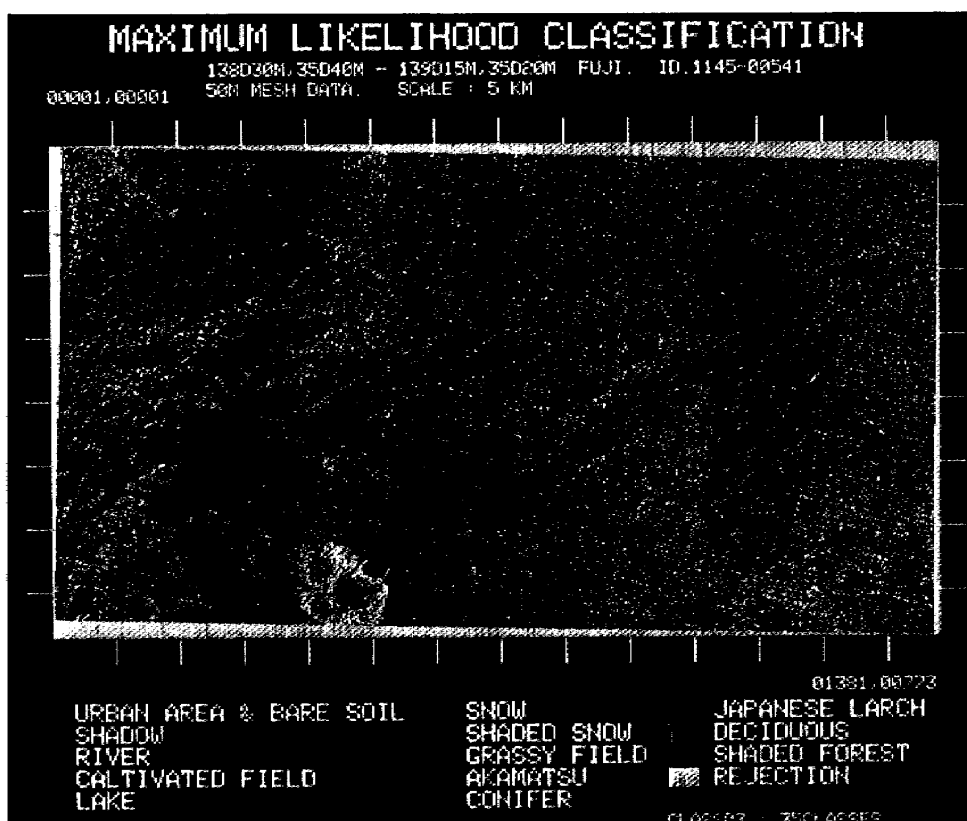


図 5-14 最尤法解析による植生分類図





表 5-12(a) 谷間の閾値  $N$  によるクラスタ構成の変化 (全体画像からのクラスタ)

$N = 1$	→	$N = 5$	→	$N = 50$
[ 草 - 1.47, 451539 ]				
[ 草 - 0.86, 155409 ]	→	[ 草 - 1.18, 602048 ]		
		[ 草 - 0.85, 284605 ]	→	[ 草 - 0.78, 889352 ]
[ 日 陰 - 0.73, 29818 ]				
[ 日 陰 - 1.32, 11148 ]	→	[ 日 陰 - 0.79, 40703 ]	→	[ 日 陰 - 0.79, 42999 ]
[ 草 - 0.84, 123377 ]				
[ 日陰林 - 1.54, 38912 ]	→	[ 草 - 0.84, 160829 ]		
		[ 草 - 1.07, 167696 ]		
		[ 草 - 0.72, 111962 ]		
		[ 草 - 0.94, 89343 ]		
		[ 草 - 0.55, 57635 ]		
		[ 草 - 1.60, 75163 ]		
		[ 草 - 0.87, 85777 ]		
		[ 耕 地 - 1.19, 94353 ]		
		[ 耕 地 - 1.00, 28954 ]		
		[ カラマツ - 0.58, 211173 ]		
		[ カラマツ - 1.41, 59724 ]		
		[ カラマツ - 1.48, 10946 ]		
		[ 日 陰 林 - 1.31, 23044 ]	→	[ 草 - 0.34, 1176599 ]
		[ 日 陰 雪 - 1.04, 86030 ]		
		[ 日 陰 雪 - 1.45, 145102 ]		
		[ 日 陰 雪 - 0.59, 82020 ]		
		[ 日 陰 - 0.96, 129552 ]		
		[ 湖 - 0.67, 103195 ]		
		[ 溶 岩 地 - 0.76, 129528 ]	→	[ 日陰雪 - 0.69, 675427 ]
		[ 広 葉 樹 - 1.06, 110857 ]		
		[ 広 葉 樹 - 1.31, 263584 ]		
		[ 草 - 1.92, 88118 ]	→	[ 広葉樹 - 0.69, 467523 ]
		[ 針 葉 樹 - 1.51, 86907 ]		
		[ 日 陰 林 - 2.04, 97738 ]	→	[ 針葉樹 - 1.37, 185786 ]

[ クラス名-D, P ] : D ; クラスタ距離 ( Separability )

P ; 画素数

表 5-12(b) 谷間の閾値  $N$  によるクラス構成の変化 (補正部分画像からのクラス)

$N = 1$	$\rightarrow N = 5$	$\rightarrow N = 10$	$N = 50$
[ 草 -1.27, 5943]			
[ 草 -1.13, 3977]			
[ 草 -1.30, 2797] $\rightarrow$ [ 草 -0.61, 12388] $\rightarrow$ [ 草 -0.62, 12471]			
		[ 草 -1.39, 6362] $\rightarrow$ [ 草 -0.61, 21338]	
[ 広葉樹 -0.59, 7528]			
[ 広葉樹 -1.55, 2457] $\rightarrow$ [ 広葉樹 -0.58, 14016]			
	[ 耕地 -0.75, 9709] $\rightarrow$ [ 広葉樹 -0.48, 23952]		
		[ 広葉樹 -1.64, 24316] $\rightarrow$ [ 広葉樹 -0.61, 52772]	
	[ 針葉樹 -1.74, 8377]		
	[ 広葉樹 -1.74, 12933] $\rightarrow$ [ 針葉樹 -1.95, 21310]		
		[ 草 -1.25, 25557] $\rightarrow$ [ 草 -1.53, 46867]	
[ ウラジ ロモミ -2.50, 6860] $\rightarrow$ [ 耕地 -2.50, 6998] $\rightarrow$ [ 広葉樹 -2.35, 8919] $\rightarrow$ [ 広葉樹 -1.54, 16026]			
		[ 日陰林 -1.56, 28168]	
		[ 日陰林 -0.64, 24856]	
		[ 日陰林 -1.12, 8248]	
		[ 日 陰 -0.19, 17355]	
		[ 日 陰 -1.07, 15365]	
		[ 湖 -0.97, 15368]	
		[ 湖 -0.52, 29219]	
		[ 日陰雪 -0.89, 57423] $\rightarrow$ [ 日陰林 -0.92, 236546]	
		[ 日陰林 -0.56, 19837]	
		[ アカマツ -0.43, 8737]	
	[ ウラジロモミ -0.45, 13014] $\rightarrow$ [ アカマツ -0.11, 35946]		
	[ 草 -0.84, 70008]		
	[ 草 -0.79, 26080]		
	[ 草 -0.93, 11374]		
	[ 草 -0.86, 8938]		
	[ カラマツ -0.81, 19801] $\rightarrow$ [ 草 -0.51, 129672]		
	[ カラマツ -1.23, 25557]		
	[ カラマツ -0.42, 19751]		
	[ アカマツ -0.93, 6534] $\rightarrow$ [ カラマツ -0.45, 27999]		
	[ カラマツ -1.64, 7940]		
	[ カラマツ -1.66, 4119]		
	[ アカマツ -0.59, 7339] $\rightarrow$ [ カラマツ -1.60, 12059]		
	[ 草 -1.91, 14331]		
	[ 日陰林 -1.85, 13600] $\rightarrow$ [ 日陰林 -1.88, 28189]		

[ クラス名 - D, P ] : D ; クラス距離 (Separability)

P ; 画素数

## 6. 結 論

従来、衛星画像はその大規模性のために、クラスタ解析にあたって一定のサンプリング処理が不可欠であったり、用い得る手法も限られるなど不満足にしか行なえなかった。これに対して、筆者らは衛星画像の大規模性を逆に有効利用し、いかなるサンプリング処理も必要ないヒストグラム・モード法を提案した。

ヒストグラム・モード法の基本的手法は、大規模データの多次元ヒストグラムの極大点探索、及び極大点から周囲へのクラスタ領域成長による谷間の決定からなる。ヒストグラム・モード法は確率密度関数のモードを把握する階層的モード法を以下の点で発展させたものである。

- (1) 大規模データのモード解析を可能にした。
- (2) 不明確な概念“データ密度”に代わり、確率密度の推定値であるヒストグラムの頻度値を用いる。
- (3)  $N$ 個のデータに対し  $N^2$  のオーダーを要していた計算負荷を  $N$  のオーダーに軽減した。

ヒストグラム・モード法の解析能力を示すために、2次元、4次元の疑似データ及びランドサットMSデータの解析実験を行なった。疑似データの解析では4次元データにおいても十分な解析精度を持つことが示された。又、莫大なランドサットデータの分類解析に十分適用可能であることが示された。

更に、ヒストグラム・モード法によるクラスタ解析の基礎となる多次元ヒストグラムの効率的構成方法に関して、従来用いられていたハッシング法を改良した多層ハッシング法を提案した。多層ハッシング法は複数のヒストグラム表を動的に使用し、全体としての処理効率、ヒストグラムの忠実度を向上させるものである。衛星画像のように、各クラスタの空間分布に偏りがあるデータのヒストグラム構成には、極めて有効である事が示された。

## 謝 辞

本研究の一部は科学技術庁特別研究促進調整費によって行なったものである。本研究を行なうにあたり計算センターの本間幸造技官及び㈱センチュリー・リサーチ・センターの田中正明、黒岩好夫氏に協

力して頂いた。ここに記して感謝の意を表わします。

## 参 考 文 献

- 1) 奥野忠一他, 「多変量解析法」日科技連.
- 2) M. Goldberg & S. Shlien, "A Clustering Scheme for Multispectral Images", IEEE, SMC-8, No.2, pp. 86—92, 1978.
- 3) S. Shlien, "Practical Aspects Related to Automated Classification of LANDSAT Imagery Using Lookup Table" CCRS TR-75-2, 1975.
- 4) K.S. Fu et al., "Information Processing of Remotely Sensed Agricultural Data" Proc. of IEEE vol. 37, No.4, pp. 639—653, 1969.
- 5) E.P. Kan et al., "The JSC Clustering Program ISOCLS and its Application" Proc. Machine Processing of Remotely Sensed Data, Purdue Univ., 1973.
- 6) R.O. Duda & P.E. Hart, "Pattern Classification and Scene Analysis", John Wiley & Sons, 1973.
- 7) J.C. Simon, "Some Current Topics in Clustering in Relation with Pattern Recognition" Proc. of 4th IJCPR, pp. 19—29, 1978.
- 8) K. Matsumoto M. Naka & H. Yamamoto, "A New Clustering Scheme for LANDSAT Images Using Local Maximum of a Multi-Dimensional Histogram" Proc. of 7th Machine Processing of Remotely Sensed Data, Purdue Univ., 1981.
- 9) 津田孝夫, 「モンテカルロ法とシミュレーション」, 培風館.
- 10) 地球観測画像情報技術データブック, 宇宙開発事業団地球観測センター, 1981.
- 11) P.H. Swain & S.M. Davis edited, "Remote Sensing The Quantitative Approach", McGraw Hill, 1978.
- 12) G.H. Ball & D.J. Hall, "ISODATA, A Novel Method of Data Analysis and Pattern Classification", SRI Tech. Report, April, 1965.
- 13) D. Wishart, "Numerical Classification Meth-

- od for Deriving Natural Classes”, Nature, vol.221, Jan. 1969, pp. 97–98.
- 14) D. Wishart, “An Algorithm for Hierarchical Classifications”, Biometrics, vol.25, 1969, pp. 165–170.
- 15) 和達清夫他, 「リモートセンシング」, 朝倉書店, 1976.
- 16) 松本甲太郎, 中正夫, 山本浩通, 「衛星画像用大規模データのクラスタリング手法」, 第24回航空宇宙技術講演会, 1980.
- 17) 松本甲太郎, 中正夫, 山本浩通, 「衛星画像クラスタリング・プログラム CASIM(I),(II)」, 第6回リモートセンシング・シンポジウム, 1980.
- 18) 松本甲太郎, 中正夫, 山本浩通, 「大規模画像の多次元ヒストグラム作成手法—多層ハッシング法—」, 第7回リモートセンシング・シンポジウム, 1981.
- 19) 松本甲太郎, 中正夫, 山本浩通, 「大規模多次元データのヒストグラム・クラスタ解析手法」, 第24回情報処理学会全国大会, 1982.
- 20) 「LARSYS 画像処理プログラム 概説書」日本IBM, 1976.

付録(1) 2.1.2節の式(2・8)  $\mathbf{A}_i$  と  $\mathbf{A}_T$  に関して2つの2次元正規分布があるとする。

$$\text{分布 } f_1: \mu_1 = (a, 0), \quad K_1 = I$$

$$\text{分布 } f_2: \mu_2 = (-a, 0), \quad K_2 = I$$

$$P(\omega_1) = P(\omega_2)$$

このとき, 分布  $f_1, f_2$  が混合した混合分布  $f_0$  は,

$$\text{分布 } f_0: \mu_0 = (\mu_1 + \mu_2) / 2 = (0, 0)$$

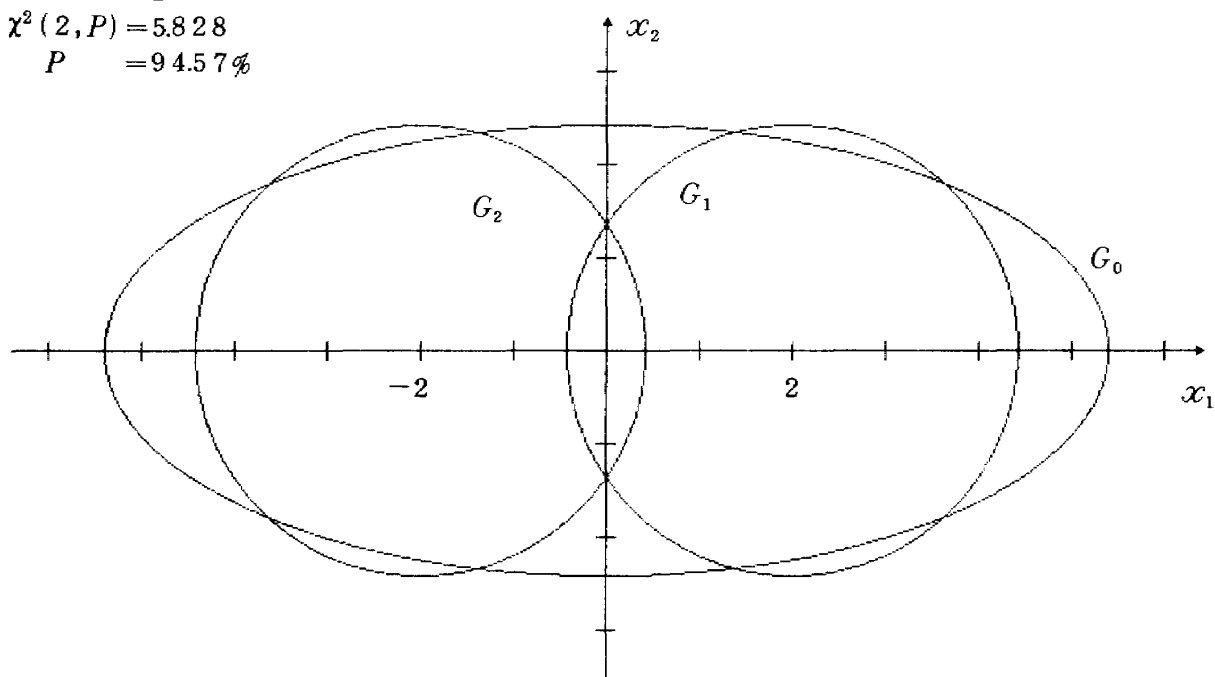
$$K_0 = \frac{1}{2}(K_1 + K_2) + \frac{1}{4}(\mu_1 - \mu_2) \cdot (\mu_1 - \mu_2) \\ = \begin{bmatrix} 1+a^2 & 0 \\ 0 & 1 \end{bmatrix}$$

分布  $f_0$  の中心部分  $P\%$  のヒストグラムを構成するとする。 $P\%$  の部分の境界は等確率密度曲線になるように定める。このとき, 各分布の  $P\%$  の領域  $G_i$  は

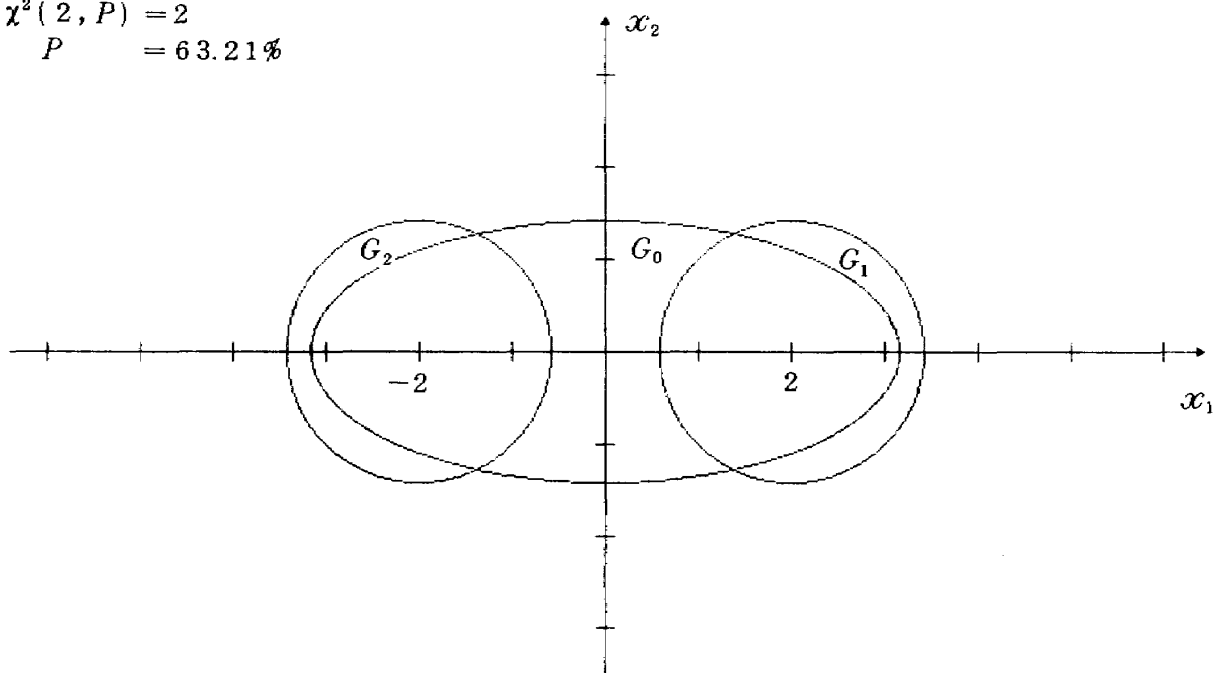
$$G_i = \{X \mid (X - \mu_i)^t K_i^{-1} (X - \mu_i) \leq \chi^2(2, P)\}$$

で定まる。 $G_0$  と  $G_1, G_2$  の包含関係は  $P$  及び両分布の距離  $a$  によって変わるが, 付図1にいくつかの例を示すように,  $P$  が大きくなると  $G_0$  は  $G_1, G_2$  を近似的に包含すると言えるであろう。

$$\begin{aligned} a &= 2 \\ \chi^2(2, P) &= 5.828 \\ P &= 94.57\% \end{aligned}$$



$$\begin{aligned} a &= 2 \\ \chi^2(2, P) &= 2 \\ P &= 63.21\% \end{aligned}$$



付図 1 混合分布と元の分布における中心部分 P %領域の違い

```

C          INITIALIZE BETA COEFFICIENTS

          CALL HBETA(IPR1,IPR2,LVEC) ← (4.14)~(4.16)

C          ***** GRAND LOOP *****

C          GET INPUT

5000 CONTINUE
          CALL SMPGT1(KF,DATA,IFR,NDT,NSMP,NDIM,LNG,INDX,MF)
          IF(LNG.EQ.0) GO TO 2000

          INDX($IDX) = INDX($IDX)+1
          ICHN(INDX($IDX))

C          ***** GENERATE HASH VECTOR *****

90 CONTINUE ← (4.4)~(4.8)

          DO 200 IV=1,LVEC
            JVS = (IV-1)*NPACK+1
            JVE = JVS+NPACK-1
            IF(IV.EQ.LVEC) JVE = MZW

            NVEC(*,*,IV) = NVEC(*,*,JVS)
            IF(JVS.EQ.JVE) GO TO 200

            DO 100 JV=JVS+1,JVE
              NVEC(*,*,IV) = NVEC(*,*,IV)*NSHFT + NVEC(*,*,JV)
100          CONTINUE

200          CONTINUE

C          ***** CALCULATE HASH KEY "A" & "B" *****

          LA = LVEC + 1
          CALL HSVEC(NVEC,BETA,IPR1,1,LA,MX,MY,LVEC) ← (4.11)
          LB = LVEC + 2
          CALL HSVEC(NVEC,BETA,IPR2,2,LB,MX,MY,LVEC) ← (4.12)
          NVEC(*,*,LA) = NVEC(*,*,LA) +1
          NVEC(*,*,LB) = NVEC(*,*,LB) +1

C          ***** HASHING SCHEME *****

400 CONTINUE
          LC = LB+1
          NDP = 1
          NDT = MXW*MYW

          DO 1000 ILP=1,NLPH

            DO 500 IDP=1,NDT
              KEY = NVEC(IDP,1,LA)
              IF(NTBL(KEY,NHT).EQ.0) GO TO 420

C          ** CHECK CONFLICT **

          DO 410 IV=1,LVEC
            IF( NTBL(KEY,IV) .NE. NVEC(IDP,1,IV) ) GO TO 440
410          CONTINUE

```

付図2 並列ハッシング・アルゴリズム

```

C          ** VECTOR MATCHING **                                ← Step 2(b)

      NTBL(KEY,NHT) = NTBL(KEY,NHT) + 1
      GO TO 500

C          ** BLANK ENTRY **                                    ← Step 2(a)

420  CONTINUE
      DO 430 IV=1,LVEC
430  NTBL(KEY,IV) = NVEC(IDP,1,IV)
      NTBL(KEY,NHT) = 1
      NCRTB = NCRTB+1
      GO TO 500

C          ** CONFLICT **                                       ← Step 2(c)

440  CONTINUE
      NVEC(NDP,1,LC) = IDP
      NDP = NDP+1

C*****NEXTKEYCALCULATION *****

500  CONTINUE
      NDP = NDP-1
      IF(NDP.EQ.0) GO TO 5000
      .....
      DO 510 I=1,LB
      510  NVEC($JDX,1,I)
      = IGAT( NVEC($JDX,1,LC) , NVEC($IDX,1,I) )
      .....
      NVEC($JDX,1,LA)
      = MOD( NVEC($JDX,1,LA)+NVEC($JDX,1,LB) , IPR1 )+1 (4.13)
      NDT = NDP
      NDP = 1

1000 CONTINUE

C          GROUND LOOP END

      GO TO 5000

      SUBROUTINE HSVEC(NVEC,BETA,IPR,IBETA,LAA,MX,MY,LVEC)

      INTEGER *4 NVEC(MX,MY,1),BETA(LVEC,2)

C          ** CALCULATE HASH KEY "A" OR "B" **
C          FOR MULTI-WORD HASHING VECTOR

      LA = LAA
      LB = LA+1                                ← (4.17)

      DO 340 I=1,LVEC
340  NVEC(*,*,LA+I) = MOD( NVEC(*,*,I) , IPR )
      IF(LVEC.EQ.1) GO TO 370

      DO 350 I=2,LVEC
350  NVEC(*,*,LA+I) = MOD( NVEC(*,*,LA+I)*BETA(I,IBETA),IPR )

      DO 360 I=2,LVEC
360  NVEC(*,*,LB) = NVEC(*,*,LB) + NVEC(*,*,LA+I)
      NVEC(*,*,LB) = MOD( NVEC(*,*,LB) , IPR )

370  CONTINUE
      RETURN
      END

```



```

SUBROUTINE TBSRT1 (NTBL, IPR1, NENT, NPRAC)

DIMENSION NTBL (IPR1, NENT)
LOGICAL    LW (NPRAC)
INDEX      $IDX/IS, IE/, $JDX/JS, JE/
INTEGER    $IDX, $JDX

ALLOCATE LW
NEND= NPRAC+1
IS   = NPRAC
IE   = NPRAC
JS   = MLP = IS-1
LW(*) = .FALSE.

C.....      ** MAIN LOOP **

DO 5000 I=1,MLP

C              Find a Position to be inserted

LW($IDX) = NTBL(JS,NENT) .GT. NTBL($IDX,NENT)

IPNT      = IFBON(LW(*))
IF(IPNT .EQ. IS) GO TO 200
IF(IPNT .EQ. 0) IPNT=NEND

C              ** DO SORT BY INSERTING **

JE   = IPNT-2
IPNT = IPNT-1

C              Block Move to insert a entry

DO 100 J=1,NENT
  ISAVE      = NTBL(JS,J)
  NTBL($JDX,J) = NTBL($JDX+1,J)
  NTBL(IPNT,J) = ISAVE
100 CONTINUE

200 CONTINUE
  IS = IS-1
  JS = IS-1

5000 CONTINUE
  FREE LW
  RETURN
  END

```

付図3 並列ソーティング・アルゴリズム

```

C.....    ** CHECK CONNECTEDNESS TO SOME AREA **

400 CONTINUE

C              Check LOWER Bound

      LW(ISDX) = RSLAND(ISDX,NLOWS).LT.DATA(1,NDPNT)+2
      K = 2
      DO 410 I=NLOWS+1,NLOWE
      LW(ISDX) = LW(ISDX)
      .AND. RSLAND(ISDX,I).LT.(DATA(K,NDPNT)+2)
410 K = K+1

C              Check UPPER Bound

      K = 1
      DO 420 I=NUPPS,NUPPE
      LW(ISDX) = LW(ISDX)
      .AND. RSLAND(ISDX,I).GT.(DATA(K,NDPNT)-2)
420 K = K+1

C              Get Number of Connected Areas

      LS = IONC(LW(ISDX))
      IF(LS.EQ.0) GO TO 500
      LSTC(LDX) = IDXL(LW(ISDX))
      GO TO 600

500 CONTINUE

CASE ***      Isolated Cell

C.....
C              Make a New Area
C.....
600 CONTINUE
      IF(LS.GT.1) GO TO 700

CASE ***      Connected to only 1 Area

C.....
C              Merge the Cell to the Area
C.....

700 CONTINUE

CASE ***      Connected to Some Areas

C.....
C              Decide the most suitable Area &
C              Merge the Cell to the Area
C              Merge the Areas if possible
C.....

```

付図4 領域とセルの連結計算の並列化



---

## 航空宇宙技術研究所報告 854 号

昭和 60 年 3 月 発行

発行所 航空宇宙技術研究所  
東京都調布市深大寺東町 7 丁目 44 番地 1  
電話武蔵野三鷹(0422)47-5911(大代表)〒182  
印刷所 株式会社 実業公報社  
東京都千代田区九段南 4-2-12

---

