

自己分配方式による新クラスター分析法の考案と
その白血球自動分類への応用

学位申請者

山口大学大学院医学系研究科保健学専攻
生体情報検査学領域

佐藤正一

目次

第 1 章 序 章	2
1-1 本論文の背景	2
1-1-1 クラスタについて	2
1-1-2 臨床検査領域におけるクラスタ分析	2
1-2 本研究の目的	3
1-3 本論文の構成	3
第 2 章 クラスタ分析	6
2-1 クラスタ分析の概説	6
2-2 階層型クラスタ分析	7
2-3 非階層型クラスタ分析	8
2-3-1 K-means 法	8
2-3-2 自己収束形アルゴリズム	11
2-3-3 K-means++法	12
2-3-4 最大距離アルゴリズム	12
2-3-5 Mahalanobis 距離を用いた K-means 法	13
2-3-6 Fuzzy c-means 法	14
2-3-7 Mahalanobis 距離を使用した Fuzzy c-means 法	15
2-3-8 EM アルゴリズム	15
2-4 混合分布モデルにおけるクラスタ数の検出	19
2-5 極小領域における局所最適解の影響	20
2-6 クラスタ分析結果の解釈	21
2-7 まとめ	21
第 3 章 クラスタ解析問題の解決策	24
3-1 画像処理方式クラスタ探索法	24
3-2 自己分配方式クラスタ分析法(SPC: self-partition clustering algorithm)	27
3-2-1 SPC 法の理論	27
3-2-2 SPC の統計量計算式	27
3-3 二次元反復切断補正法 (ITC: Iterative truncation-correction method)	28
3-3-1 SPC/ITC 法の手順	29
3-3-2 補正係数の算出	30
3-4 開発したプログラムについて	32
3-5 各クラスタ分析法の比較検証法	33
3-5-1 シミュレーションデータによる IP 方式クラスタ探索法の性能評価	33

3-5-2	近接するクラスターにおける SPC/ITC 法の性能評価	34
3-5-3	近接するデータ数(密度)や広がりを変化させたクラスターに対する各種クラスター分析法の比較	36
3-5-4	EM アルゴリズムと SPC/ITC 法のバックグラウンドノイズの影響	39
3-5-5	白血球分画疑似モデルによる一致率	40
3-6	まとめ	44
第 4 章	臨床検査分野への応用	48
4-1	泳動分析	48
4-1-1	各種泳動検査の原理	48
4-1-2	泳動分析の問題点	49
4-1-3	自動血球分析装置の臨床的有用性	50
4-1-4	白血球分類・フローサイトメリーの原理	50
4-1-5	白血球分類フローサイトメリー二次元散布図の特徴	53
4-1-6	塗抹鏡頭における白血球分画の統計的限界と手技的限界	54
4-1-7	フローサイトメリーの限界	55
4-2	実データによる検証	57
4-2-1	SPC/ITC 法と EM アルゴリズムの性能を評価	57
4-2-2	健常人検体による目視法、機器分析、SPC/ITC 法の比較	59
4-2-1	異常検体による目視法と Pentra MS CPR、SPC/ITC 法の性能	60
4-2-2	Pentra MS CRP の細胞区画を超えた SPC/ITC 法の分析結果	64
4-3	白血球分画の実データによる各種クラスター解析のまとめ	66
第 5 章	総 括	69
5-1	本研究の成果	69
5-2	本研究の限界	70
5-3	本研究の今後の展開	71
5-3-1	SPC/ITC 法から得られたパラメータの利用	71
5-3-2	SPC/ITC 法の白血球分画以外の臨床分野への応用	71
5-3-3	疾患分類への応用	72
5-3-4	統計分野への応用	72
謝 辞	72
付録	76
●	Mahalanobis 距離について	76
●	大域的最適解と局所的最適解	78
●	ベイジアン情報量規準 (BIC: Bayesian Information Criterion)	79

- 目視法の結果と Pentra MS CPR の測定データおよび SPC/ITC 法適応図.....80

図目次

図 2-1	クラスター分析の階層的手法と非階層的手法の例	7
図 2-2	群平均法の概念図.....	8
図 2-3	二次元散布図による K-means 法の解析アルゴリズム計算過程図	10
図 2-4	K-means 法の初期値の違いによる最終結果への影響.....	11
図 2-5	最大距離アルゴリズムの初期中心点の設定概念図.....	12
図 2-6	K-means 法のユークリッド距離と Mahalanobis 距離(MD)の対比.....	14
図 2-7	K-means 法と Mahalanobis 距離を用いた K-means 法および EM アルゴリズムのクラ スター分析の関係	16
図 2-8	混合正規分布の例.....	18
図 2-9	EM アルゴリズムの極小領域での局所最適解と	18
図 2-10	EM アルゴリズムの極小領域での局所最適解と	19
図 2-11	クラスターが明確な場合の BIC によるクラスター数の検出例.....	20
図 2-12	極小領域における局所最適解の影響.....	21
図 3-1	IP 方式クラスター探索法の初期クラスターを検出するための.....	26
図 3-2	二次元反復切断補正法理論の概念図	29
図 3-3	反復切断補正法の算出方法	30
図 3-4	二次元反復切断補正法の補正係数.....	31
図 3-5	SPC/ITC 法の概念図.....	31
図 3-6	開発したプログラムの実行結果	32
図 3-7	初期クラスターを識別するための IP 方式クラスター探索法の性能.....	34
図 3-8	SPC/ITC 法のクラスターのオーバーラップの影響	35
図 3-9	EM アルゴリズムによるクラスターのオーバーラップの影響.....	35
図 3-10	近接するクラスターにおける各種クラスター分析法の比較	37
図 3-11	R 言語プログラミングの mclust で分析結果	38
図 3-12	EM アルゴリズムと SPC/ITC 法のバックグラウンドノイズの影響.....	39
図 3-13	今回の検討に使用した血球自動分析装置.....	40
図 3-14	健常人データを使用して分析装置から出力された二次元散布図と人工的に作成し た白血球分画疑似モデル.....	41
図 3-15	白血球分画疑似モデルの RD 行列と IP 方式クラスター探索法.....	41
図 3-16	白血球分画疑似モデルにおける各種クラスター分析法の検証結果.....	42
図 4-1	血清タンパク分画のデンシトグラムと各分画に属する主要血清タンパク成分.....	48

図 4-2	血清タンパク分画像	49
図 4-3	1次元混合分布における境界値の問題点	49
図 4-4	フローサイトメータの模式図	51
図 4-5	フローサイトメトリーの光学的データ処理図	51
図 4-6	堀場製作所の自動血球分析装置 Pentra MS CPR からの二次元散布図	52
図 4-7	シスメックス社の自動血球分析装置 XN-1000 からの二次元散布図	52
図 4-8	シスメックス社の XN-1000 の各細胞分布と細胞形態	53
図 4-9	細胞比率 5% の 200 カウントと 400 カウントの 95% 信頼区間の比較	54
図 4-10	ウェッジ法による血液塗抹標本の状態	54
図 4-11	XN-1000 における特徴的なデータの二次元散布図	56
図 4-12	100 人の無作為抽出した各細胞のクラスター重心点をプロットした図	57
図 4-13	Pentra MS CRP における SPC/ITC 法と EM アルゴリズムの性能評価	58
図 4-14	XN-1000 における SPC/ITC 法と EM アルゴリズムの性能評価	59
図 4-15	健常検体による目視法、Pentra MS CPR、SPC/ITC 法の比較	62
図 4-16	異常検体による目視法、機器分析、SPC/ITC 法の比較	63
図 4-17	Pentra MS CRP の細胞区画を超えた時の SPC/ITC 法のクラスター分析結果	65

表目次

表 2-1	階層型クラスター分析距離の説明	8
表 2-2	BIC を求めるための分散モデル	20
表 3-1	相対密度作成に使用したフィルタ行列 F	25
表 3-2	白血球分画疑似モデルによる各クラスター分析法の誤差率	43

第 1 章

序 章

- 1-1 本論文の背景
- 1-2 本研究の目的
- 1-3 本論文の構成

第 1 章 序 章

1-1 本論文の背景

1-1-1 クラスタについて

クラスター(cluster)とは、物質科学では原子や分子の数個から数十個からなる集合を意味し、天文学では恒星の集団、銀河の集団という意味であり、コンピュータの領域ではオペレーションシステムの管理する最小単位など様々な意味がある。これから述べるクラスターとは、“集団”という意味である。

クラスター分析は、外的な基準なしで、観測されたデータの類似性をもとに似たもの同士をまとめるデータ分析手法で、データ解析における自動分類手法である。クラスター分析は 20 世紀後半から盛んに研究され、多くの応用技術が存在し、人類学、生物学、医学、心理学、統計、数学、工学、コンピュータサイエンスなどの幅広い科学分野で利用されている。

クラスター分析の基本的アルゴリズムは、階層的クラスター分析と非階層的クラスター分析(分割最適化手法)に分けられる。階層的クラスター分析は、各点の距離について樹形図(デンドログラム)を使ってクラスター分析を行うことから階層的的手法と呼ばれる。階層的クラスター分析は、全サンプルのペアの間の距離を計算するため、計算量はサンプル数の二乗に比例することになり、多くのデータを処理するために時間を要するという欠点があるが、生物学、遺伝子学、形態学等の分野で利用されている。一方、非階層的クラスター分析は、乱数を使ってランダムに仮の重心点を設定して、その重心点に近いデータを集めて一旦仮のクラスターとし、クラスターの重心点の再計算を繰り返すものである。階層的クラスター分析に比べ非常に高速に行うことが可能な手法で、画像処理等の高速な処理が必要な分野で利用されている。

1-1-2 臨床検査領域におけるクラスター分析

タンパク電気泳動分析やフローサイトメトリー法を使った白血球分類など、パターン認識が要求される検査において、クラスター分析は必要性の高い技術と言える(図 1-1)。しかし、臨床検査領域のデータは、多種多様な異常データの混入が発生することから、現状のクラスター分析では対応ができない。一方、タンパク成分や血球成分は、その分析技術については非常に高度な測定技術を用いて解析可能な状況となっているにも関わらず、的確なクラスター分析を行っていないために、誤分類を生じやすいという問題を残したままとなっており、臨床分野で使用できるクラスター分析手法の開発が望まれる。

例えば、タンパク分画では、様々な蛋白質が混合した混合分布データであると考えられ、その画像も明確なデータの切れ目がなく連続したものである(図 1-1 上段)。この混合分布に対して谷値で分画化し、各種蛋白成分(アルブミン分画、 α_1 分画、 α_2 分画、 β 分画、 γ 分画)のパーセント(%)を算出する手法が一般に利用されているが、混合分布ではクラスターの分布幅を考慮して対処する必要性がある。

1-2 本研究の目的

本研究は、クラスター分析において、初期クラスターの推定とクラスター境界領域の処理を的確に行うことで、分類結果の最適化を行うことにある。特に、現行のクラスター分析では、臨床検査領域で発生する極端なデータ数の違い、各種のデータノイズ(細胞分類では細胞の断片や凝集によって発生する)、高密度データなどに対応できない。これらの問題点を解決することで、臨床検査領域においてクラスター分析の応用を可能とし、分類や分画化を必要とする臨床検査の自動化に貢献することができる。具体的な臨床検査分野への応用例としては、自動血球分析装置における白血球分類があり、従来は血球の二次元散布図の固定領域に入った細胞数をそのまま用いて分類しているため、各種病態による変化への対応が不十分であった。しかし、新しく開発したクラスター分析法を適用すると、それらの問題が明快に解消できることをシミュレーションや実臨床データ事例を紹介する。

なお本研究は、臨床検査領域ばかりでなく、統計、工学、コンピュータサイエンスなど応用範囲の広いクラスター分析法を目指している。

1-3 本論文の構成

本稿は以下のように構成されている。第 2 章では、基本的なクラスター分析手法である階層型クラスター分析と非階層型クラスター分析について概説し、さらに個別のクラスター分析法の特徴および利点・欠点について述べる。第 3 章では、新開発した3つのクラスター分析法のアルゴリズムについて解説を行う。具体的には、①非階層型クラスター分析において、これまで大きな課題であったクラスター数とその初期位置の決定法に対して、それを解決するための画像処理技術であるイメージ プロセッシング方式クラスター探索法(Image-Processing (IP) algorithm for cluster search)について、②クラスター間の重なりやノイズの影響を軽減するための二次元反復切断補正法(ITC: iterative truncation-correction)と自己分配方式クラスター分析法(SPC: self-partitioning clustering)の解説である。第 4 章では、本研究の臨床検査データへの応用の具体例として、フローサイトメトリー法による白血球自動分類を例に挙げ、分析装置よ

り得られた二次元散布図の実データを、各種クラスター分析法で分類し、その分析結果を方法間で比較する。第 5 章は、総括と今後の展開として、本クラスター分析法の病態解析への応用に関して述べる。

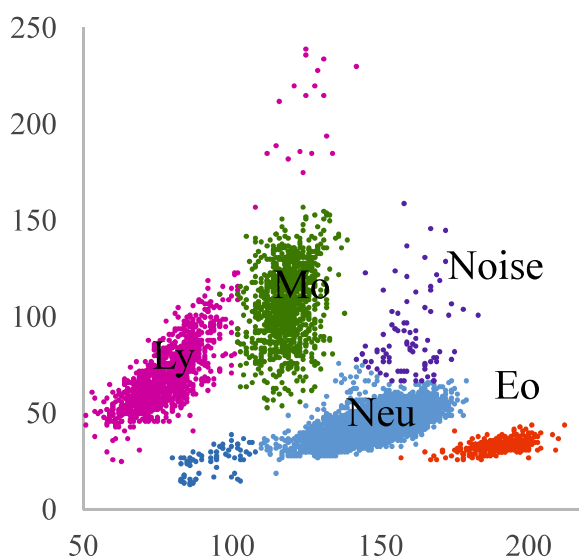
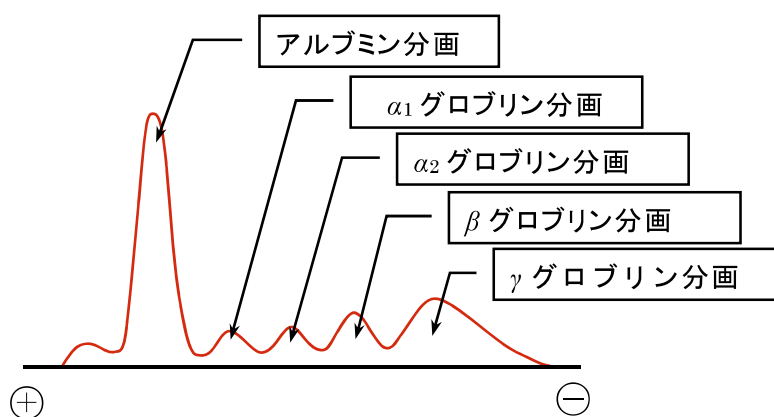


図 1-1 臨床検査においてクラスター分析が要求される分野
血清タンパク分画像(上)と自動血球分析装置の白血球分類(リンパ球[Ly], 単球[Mo], 好中球[Neu], 好酸球[Eo], ノイズ[Noise])の二次元散布図(下)

第 2 章

クラスター分析

- 2-1 クラスター分析法の概説
- 2-2 階層型クラスター分析
- 2-3 非階層型クラスター分析
- 2-4 混合分布モデルにおけるクラスター数の検出
- 2-5 クラスター分析結果の解釈
- 2-6 まとめ

第 2 章 クラスター分析

本章では、各種のクラスター分析法について解説する。

2-1 クラスター分析の概説

クラスターへの分類を自動的に行うには、何らかの仮定を導入する必要があり、一般的には、各データ間の距離が小さいもの同士をまとめるようにクラスターへの分類が行われる。クラスター分析法には、大きく分けて階層型クラスター分析と非階層型クラスター分析がある。

階層型クラスター分析は、最も類似した(距離が近い)データを選び、それを逐次的にクラスターへまとめていくことで階層的なクラスター構成を得る手法である。最終的には図 2-1 に示すような樹形図(デンドログラム)としてまとめられる。階層型クラスター分析の特徴としては、関係性の近いものから順番にくっつけていく方法であるため、あらかじめクラスター数を決めておく必要がないことが上げられる。図 2-1 の左パネルに示すように、どのレベルで分類するかを後から適時選択することが可能である。

一方、非階層型クラスター分析は、あらかじめいくつのクラスターに分けるかを決めておき、データをそのいずれかに所属させる形で分類していく方法である。巨大なデータであっても、高速に処理できるという特徴を持っている。なお、非階層的クラスター分析は、人工知能における機械学習において、外的基準のない条件で分類する「教師なし学習」に相当するもので、事前にどのクラスターに属するかを規定せず、データの類似性または相違性に基づいて、各クラスターの領域とそれに所属するデータを、動的・試行錯誤的に、反復処理により最適化するものである¹⁾。

代表的な非階層的クラスター分析である K-means 法は、1957年に E. Forgy が発表し、J. MacQueen らが命名した方法である^{2,3)}。評価関数を定めて、データの分割が最適となる複数のクラスター重心点を探索するもので、アルゴリズムは極めて簡単なものである。その他、K-means 法から派生した Fuzzy c-means 法や EM アルゴリズムがある。本章では、各種のクラスター分析法について紹介する。

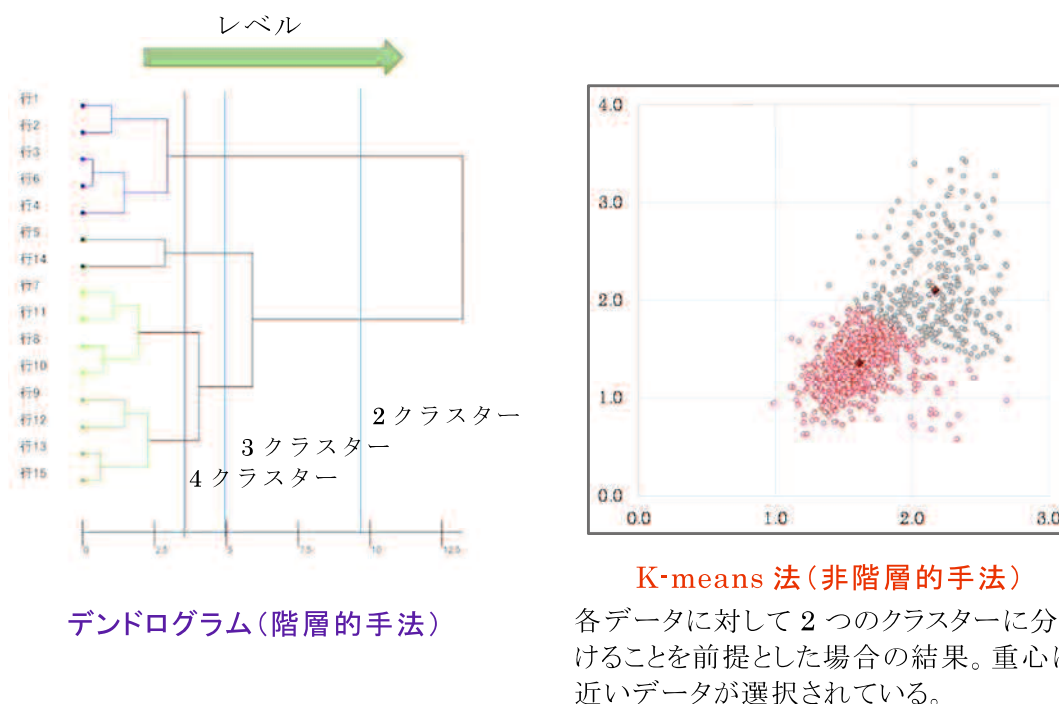


図 2-1 クラスター分析の階層的手法と非階層的手法の例

2-2 階層型クラスター分析

階層型クラスター分析は、クラスター間の距離関数に基づいて、最も距離の近い二つのクラスターを逐次的に併合することによってクラスターに分けるもので、全ての対象が一つのクラスターに併合されるまで繰り返す手法である。データ数が多い場合には計算時間が膨大となることやデンドログラムが巨大となり結果が不明瞭となり、デンドログラムを作成することが困難な場合がある。

距離関数には、最短距離法、最長距離法、群平均法、ウォード法(クラスター距離の二乗の総和を最小化する)等がある。各距離関数の特徴は、最短距離法では空間濃縮(一旦できたクラスターの影響を大きく受けて偏ったクラスター形成となってしまう現象)という性質のため、極めて外れ値に弱く、実データでは良い結果が得られない。逆に、最長距離法には併合されてできたクラスターが、以降のクラスター形成において併合されにくくなるという空間拡散という性質があるため、極端値が別クラスターとして扱われやすくなる。その他の方法は極端なクラスターが形成されやすい傾向がある。一般的な階層型クラスター分析の距離計算法として、ウォード法が最も分類感度が高いとされる。表 2-1 に各距離計算法を示す。

表 2-1 階層型クラスター分析距離の説明

手 法	特 徴
単連結法 (最短距離法)	クラスター間で各クラスターに属する対象物間の距離が最短のものをそのクラスターの距離とする。
完全連結法 (最長距離法)	クラスター間で各クラスターに属する対象物間の距離が最長のものをそのクラスターの距離とする。
群平均法	非加重結合法とも言う。各クラスターに属する全対象物間の距離の平均をクラスター間の距離とする。
重心法	クラスター間の重心間距離からクラスター間距離を導く。別のクラスターとの距離は重心間距離を要素数で重み付けした点から生成する。
メディアン法	重心法と若干異なり、重心間距離を要素数で単純に中点から別のクラスターとの距離を生成する。
Ward 法 (最小分散法)	クラスター内の距離平方和の増加量が最小になるクラスターを併合させる。最も明確なクラスターを作り、分類感度が高い。

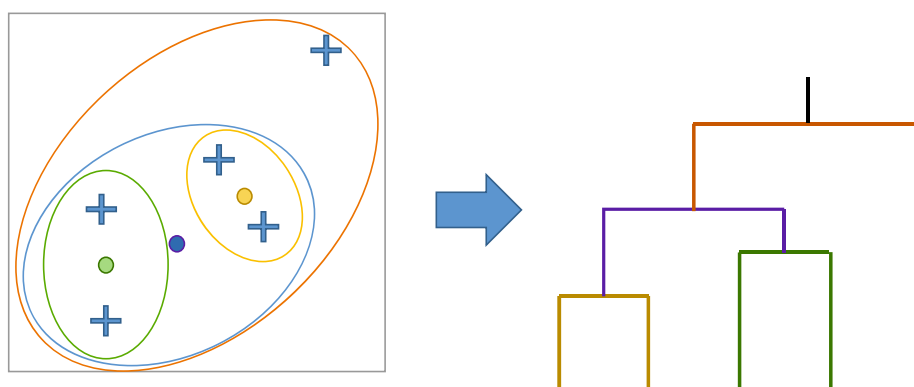


図 2-2 群平均法の概念図

データポイントの近いものを見つけてクラスターを作り、順次拡大して、距離に応じて樹形図(デンドログラム)として表現する。

2-3 非階層型クラスター分析

非階層型クラスター分析では、通常データ数 n に対して解の候補は指数的に増加するので、一般に大域最適解は計算できない。代わりに局所最適解を求めて近似する方法である(大域最適解と局所最適解については付録に記す)。

2-3-1 K-means 法

非階層型クラスター分析の代表的な方法として E. Forgy が発表した K-means 法が有名である²⁾。K-means 法は、各データは必ず距離が最小のいずれか一つのクラスターに所属

するという条件のもと、反復計算で各クラスターの所属データとその重心を最適化していく手法である。手順は、 K 個の初期クラスター位置(重心)をランダムに決めて、分類すべきデータを、最も距離が近いクラスターに所属させる。次に、各クラスターの重心を、それに所属するデータに基づいて再計算で修正し、これに応じて全データの所属を更新する。この処理を反復すると、次第に各クラスターの重心は、データに対して最適化されるというものである。この単純なアルゴリズムにより、データを自動分類する。ここで、距離指標に何も用いるかが重要で、それによって結果が変わる。

K-means 法は、式(2.1) の評価関数 φ を最小化するクラスター重心を見つけることで、データをあらかじめ設定した K 個のクラスターに分割するものである。理想的クラスターとしては、各クラスターの周辺データがオーバーラップしないことが望まれる。

$$\varphi = \sum_{x_j \in X} \min_{i \in k} \|x_j - c_i\|^2 \quad (2.1)$$

x_j は各データを示し、 n はデータの総数、 c_i は各クラスター重心を示す。

K-means 法の解析アルゴリズムは、以下の方法により行い、図 2-3 に計算過程を図示する。

- 1) クラスター数 K を決める。
- 2) 分布の重心の初期値をランダムに設定する。
- 3) 重心からの距離(ユークリッド距離)が近いデータをその重心のクラスターに割り当てる。
- 4) 割り当てられたデータから、平均値を算出し、再計算された平均値を重心に移動し、全データの所属を修正する。
- 5) 3)~4)の計算をクラスターの評価値、平均一データ間の残差平方和 (RSS: residual sum of square)で行い、RSSの減少値が閾値より小さい、または変化しなくなるまで繰り返す。

なお、ユークリッド距離とは、ピタゴラスの公式によって与えられるもので、2次元の場合には公式(2.2)で求められ、RSS(2.3)は各ベクトルの重心からの距離の2乗をすべて合計したもので、この値が最小値になっていることでクラスターへの適合度を評価する。

$$\text{ユークリッド距離} = \sqrt{(x_i - x_0)^2 + (y_i - y_0)^2} \quad (2.2)$$

$$\text{RSS} = \sum_{g=1}^K \sum_{i=1}^n |\vec{x}_i - \vec{u}|_g^2 \quad (2.3)$$

K-means 法の利点として、アルゴリズムが単純であり、非常に高速にデータ処理が行えるため、大規模データに適用可能であることから多くの分野で利用されている。欠点として、ユークリッド距離を使用した K-means 法は、バラツキが同じ大きさのクラスターに分割する性質があるため、クラスター内分散が異なるクラスターには適用が困難である。また、ランダムに設定される初期値に最終結果が大きく影響されることが多く、特にクラスターの大きさが異なる場合や、分布形状が楕円状である場合には適用できない(図 2-4)。

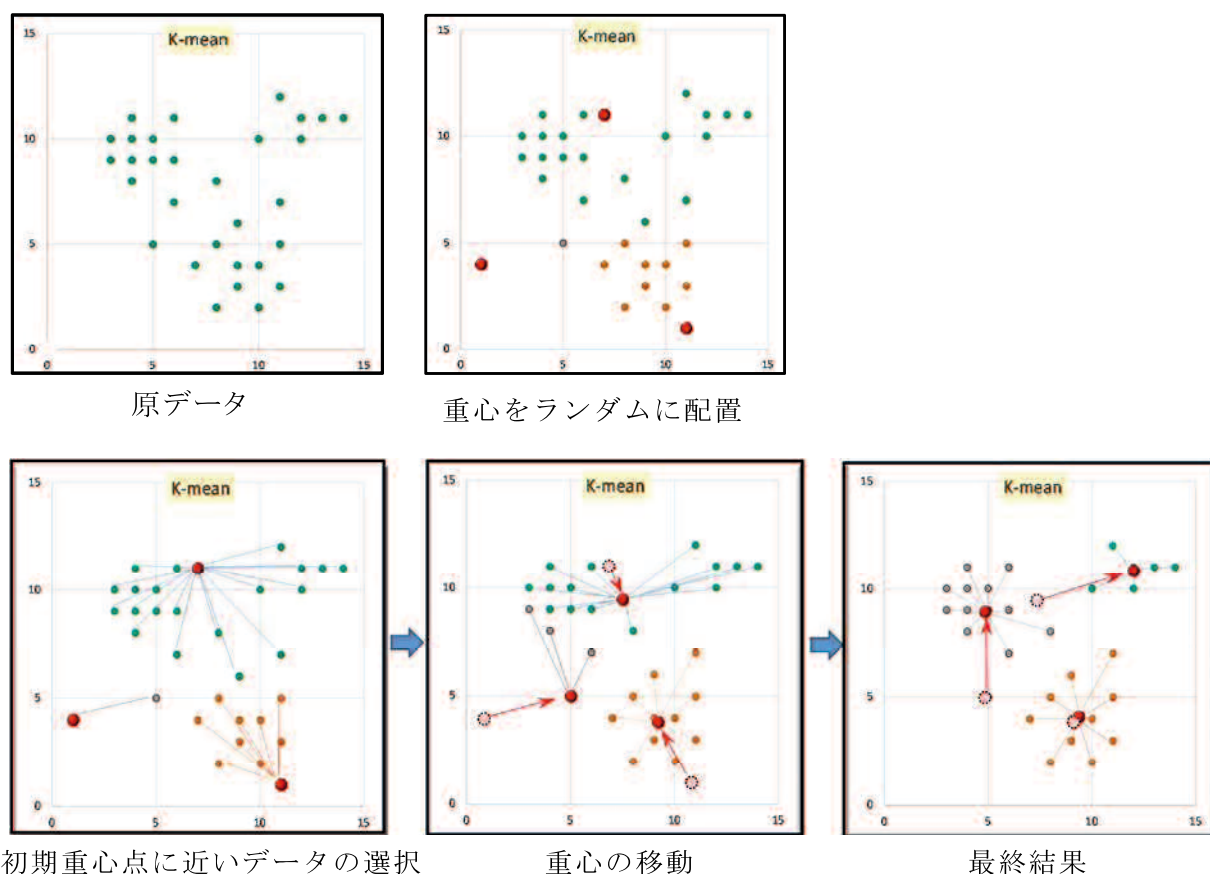


図 2-3 二次元散布図による K-means 法の解析アルゴリズム計算過程図

- 1) クラスタ数を任意に決定する (ここではクラスタ数を 3 とする)。
- 2) 乱数を使って散布図の中の 3 点を選択して、初期重心点とする(赤丸)。
- 3) 初期重心点からの距離 (ユークリッド距離)をもとに全要素を 3 群のいずれかに振り分け、振り分けられたデータを基に各クラスターの平均値を再計算する。
- 4) 平均値に基づき重心点を移動して、移動した重心点からの各点の距離を再計算する。
- 5) RSS が十分小さい、または重心点が動かなくなるまで 3)から 4)を繰り返す。

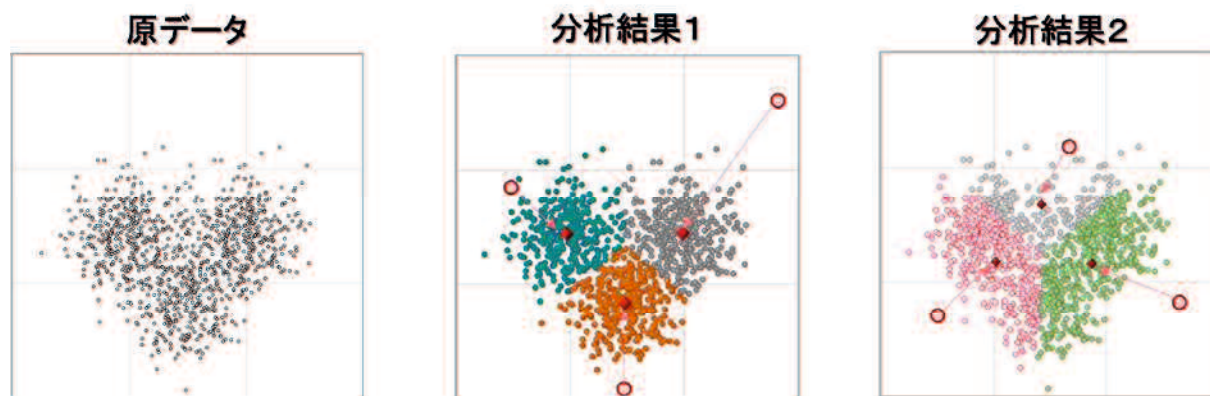


図 2-4 K-means 法の初期値の違いによる最終結果への影響

全く同じデータで、初期重心点の設定値位置によって最終重心点(赤丸)が異なった例

2-3-2 自己収束形アルゴリズム

K-means 法は初期値をランダムに設定する方法で行っているが、K-means 法の結果は局所最適解のため初期クラスターに依存する。この欠点を補う方法として、自己収束形 (ISODATA: Interactive Self-Organization of Data) アルゴリズムが考えられた⁴⁾。ISODATA アルゴリズムは、①同じクラスターに属するサンプルが閾値未満の場合、そのクラスターを作らない、②クラスター間距離が閾値未満の場合、クラスターをまとめる、③クラスター内の分散が大きくなりすぎるとクラスターを分割するという条件が加えられたものである。ISODATA のアルゴリズムを下記に示す。

ISODATA のアルゴリズム

- 1) 初期クラスターの平均ベクトルおよび、分割パラメータ、融合パラメータ、クラスターに含まれる最小要素数を設定。
- 2) 特徴空間中で、各要素を最短距離(ユークリッド距離)にあるクラスターに所属させる。
- 3) 全要素がいずれかのクラスターに属した後、各クラスターの重心ベクトルを計算。
- 4) いずれかのクラスターが、最初に設定した分割閾値を越した場合、そのクラスターを2つに分割。
- 5) 前記、第4段階で分割が生じた場合、2)に戻る。
- 6) いずれかのクラスターの要素数が、最初に設定した最小要素数より少なければ、このクラスターを削除。
- 7) 6)段階で、削除が生じた場合、2)に戻る。
- 8) いずれかのクラスター間の距離が、最初に設定した融合パラメータの値より小さい場合、そのクラスター対を融合する。
- 9) 8)段階で融合が起きた場合、2)に戻る。
- 10) 段階 2)から 9)までを、定常状態に達するまで、あるいは最初に決められた繰り返し数に達するまで繰り返す。

2-3-3 K-means++法

K-means++法は、K-means 法の初期値の問題を改良した方法で、初期クラスターの重心点はなるべく離れていた方がよいという考えに基づいたものである⁵⁾。既に決定されたクラスター重心から、ある程度遠いデータ点をクラスター重心に決定することができ K-means 法の欠点を補うものである。はじめに 1 つ目のクラスター重心をデータの中からランダムに選択する。2 個目以降の重心は、それまでに選択済みの重心と各データの最小距離に基づく次式の重み付き確率分布を満たすものを選択する。これにより最短距離が大きなもの（最大ではない）が選択される確率が高くなる。ただし、この方法であっても最初に選択されたクラスター一点に依存する問題が残されている⁶⁾。また、初期値の検索に対する計算コストが大きい。

$$\phi(x'_i) = \frac{D(x'_i)^2}{\sum_{j=1}^n D(x_j)^2} \quad (2.3)$$

2-3-4 最大距離アルゴリズム

最も離れているデータを基準としてクラスター重心の初期値として選択する手法である。最大のデータ間距離に対して任意の制約基準 n/m を設定し、その値より大きな場合は、クラスターの初期重心点とする設定法である。

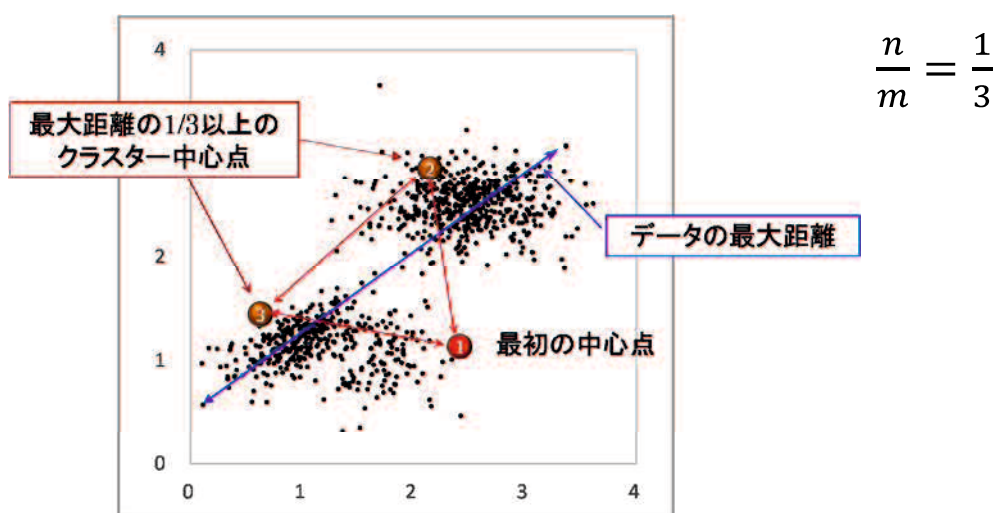


図 2-5 最大距離アルゴリズムの初期中心点の設定概念図

- 1) クラスタ数 K の数を決める。
- 2) データの最大距離を計算し、任意の制約基準を設定する。
ここでは $1/3$ とする
- 3) ランダムに最初の重心点を設定する。
- 4) 次に選択される重心点は、最大距離の $1/3$ 以上の場合にのみ選択可能とする。

2-3-5 Mahalanobis 距離を用いた K-means 法

一般的に K-means 法ではユークリッド距離を使って計算されるが、Mahalanobis 距離 (2.4)を用いた方が分類結果が良くなることが多い⁷⁾。Mahalanobis 距離は各変数間の分布の相関係数(2.6)を考量した距離で、次式によって算出される⁸⁾。

$$D = \sqrt{\frac{u_1^2 + u_2^2 - 2\rho u_1 u_2}{1 - \rho^2}} \quad (2.4)$$

$$u_1 = \frac{x_1 - \mu_1}{\sigma_1}, \quad u_2 = \frac{x_2 - \mu_2}{\sigma_2} \quad (2.5)$$

$$\rho = \frac{\sigma_{12}}{\sigma_1 \sigma_2} \quad (2.6)$$

$$\sigma_1 = \sqrt{\frac{1}{n-1} \sum n(x_{1,i} - \mu_1)^2} \quad (2.7)$$

$$\sigma_2 = \sqrt{\frac{1}{n-1} \sum n(x_{2,i} - \mu_2)^2} \quad (2.8)$$

$$\sigma_{12} = \sqrt{\frac{1}{n-1} \sum n(x_{1,i} - \mu_1)(x_{2,i} - \mu_2)} \quad (2.9)$$

Mahalanobis 距離を使用した K-means 法は、ユークリッド距離を使用した K-means 法とは大きく概念が異なる。すなわち、ユークリッド距離を使用した K-means 法では個々のデータは必ず K 個あるクラスターのいずれか一つに所属するハードクラスタリング(0 / 1 の関係)であるが、Mahalanobis 距離は所属確率(尤度)を距離とする考え方に基づいている。利点は、二次元データの場合で、2 変数間に関連性がある場合、相関性を基に距離を計算する方が単純にユークリッド距離から算出するよりもフィッティング精度が高まる点にある。図 2-6 は、大きさの異なる 2 つの楕円形のクラスターを作成し、ユークリッド距離を使用した K-means 法と Mahalanobis 距離を使用した K-means 法を比較したものである。ユークリッド距離を使用した K-means 法は、単純に重心距離が近い方にデータが所属するが、Mahalanobis 距離を使用した K-means 法は楕円に対して上手くフィットしている。Girolami は、クラスター間の分離境界が非線形である場合、K-means 法などの方法が失敗すると述べており、楕円形の分布に対してユークリッド距離は不適と言える⁹⁾。なお、Mahalanobis 距離の解説を付録に記載する。

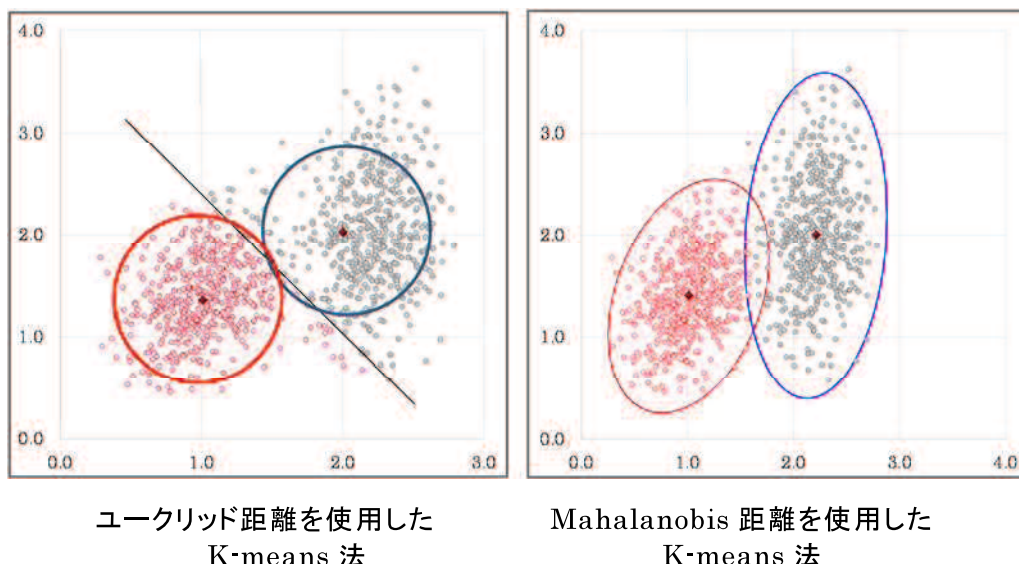


図 2-6 K-means 法のユークリッド距離と Mahalanobis 距離(MD)の対比
 各図のデータセットは、 $n = 500 : 500$ のであるが、相関関係がある楕円形の分布ではユークリッド距離よりも Mahalanobis 距離を使ったクラスター分析の方が適合度が高い。計測されたデータ数についても、K-means 法が $n = 484 : 516$ に対して、Mahalanobis 距離を使用した K-means 法では $n = 501 : 499$ で設定したクラスター数と近似していた。

2-3-6 Fuzzy c-means 法

Fuzzy c-means 法は Bezdek が報告した方法で、K-means 法がハードクラスター分析であるのに対して、Fuzzy c-means は帰属確率が $0 \sim 1$ までの間をとるソフトクラスター分析に位置づけられている手法である¹⁰⁾。すなわち、K-means 法が「所属する」か「所属しない」の二者択一の問題であるのに対して、複数のクラスターへ同時に所属することを許容する考え方である。Bezdek は K-means 法の目的関数を非線形に修正するために、Fuzzy c-means 法では所属値を任意のべき乗値 $m (m > 1)$ を導入してファジーさを与える手法を使っている。Fuzzy c-means 法は、様々な検討がなされ、現在でも派生アルゴリズムが展開されている。

Fuzzy c-means のアルゴリズムは、クラスター重心 v_i とデータとの距離は d_{ik} (2.11) で、 x_k は k 番目のデータを表す。また、パラメータ $m (m > 1)$ は曖昧さを決定するもので、 m が 1 に近づくと曖昧さは減少し、 $m=1$ のとき K-means 法の結果と同様になる。

$$J = \sum_{k=1}^n \sum_{i=1}^c (u_{ik})^m d_{ik} \quad (2.10)$$

$$d_{ik} = \|x_k - v_i\|^2 \quad (2.11)$$

$$v_i = \frac{\sum_{k=1}^n (u_{ik})^m x_k}{\sum_{k=1}^n (u_{ik})^m} \quad (2.12)$$

$$u_{ik} = \left[\sum_{k=1}^n \left(\frac{d_{ik}}{d_{jk}} \right)^{1/m-1} \right]^{-1} \quad (2.13)$$

Fuzzy c-means のアルゴリズム

- 1) クラスタ数とファジー係数 m を任意に決定する
- 2) クラスタ重心決定
K-means 法同様ランダムに、クラスタ対象点から重心を決定
- 3) クラスタ対象点ごとの、クラスタに属する割合 (u_{ik}) を更新
クラスタ対象点とクラスタ重心点ごとに帰属確率を計算
それぞれのクラスタに属する割合を決定
- 4) 各クラスタの重心点を再設定
各クラスタ対象点のクラスタ割合をもとに、新たなクラスタ重心点を設定
- 5) 2) ~ 4) を重心点が動かなくなるまで繰り返す

2-3-7 Mahalanobis 距離を使用した Fuzzy c-means 法

標準的な Fuzzy c-means 法の距離計算はユークリッド距離であるが、近年様々な距離に基づいた Fuzzy c-means 法が開発されている。Mahalanobis 距離を使用した K-means 法と同様に Fuzzy c-means 法でも Mahalanobis 距離を用いた Fuzzy c-means 法が報告されている^{11, 12)}。他に、ガウス分布の密度関数を使用するものもある¹³⁾。これらの方法は、K-means 法のような二値的で明確な判別ではなく、確率によってクラスタの属性を示すのに有効な手段である。

2-3-8 EM アルゴリズム

期待値最大化法 (EM アルゴリズム: Expectation Maximization Algorithm)

EM アルゴリズムは、観測されたデータに対してデータ構造の仮説 (正規分布など) をもとに、尤度を最大化する方法を用いてクラスタ分析を行うもので、Dempsterらが考案した方法である¹⁴⁾。EM アルゴリズムは、非線形な分布の場合に数値的な解法が出来ない問題に対して、複雑な分布を正規分布の複合体と考えて、最尤推定法を用いて逐次的に解決しようとするものである。複数の正規分布が混在する混合分布データに対して、教師なしクラスタ分析手法として有用な分析手法である。

EM アルゴリズムは、K-means では解決できないオーバーラップの問題に対して対応する手法と言える。EM アルゴリズムは正規分布を仮定した確率分布 (混合比) でクラスタを求めていくため、クラスタ間に厳密な境界が存在しない。すなわち、オーバーラップするクラスタを確率的に区分するモデルを推定方法と言える。得られたモデルは、密度推定と判別分析などの多変量解析に使用することができる。

各アルゴリズムを整理するためにユークリッド距離による K-means 法と Mahalanobis 距離を用いた K-means 法および EM アルゴリズムの関係を図 2-7 に示す。パネル A は作成した相関性のある 3 つのクラスタとその密度曲線を示し、パネル B は密度関係を表す立体

図である。パネル C は K-means 法による結果で、3 つのクラスターに対してユークリッド距離を基にクラスターが分けられている。各クラスターの距離が等距離となるように計算されるために、クラスターの広がり方が円形でない場合にはクラスター分析が失敗する。パネル D は Mahalanobis 距離を用いた K-means 法で、クラスターが楕円形の時(相関関係がある場合)に適した方法であるが、高度にオーバーラップするクラスターの場合には分析が失敗する。一方、EM アルゴリズムでは複雑に入り組んだクラスターであっても、クラスターを捉えられており、オーバーラップするクラスターに対しても上手くクラスター分析が行われる。

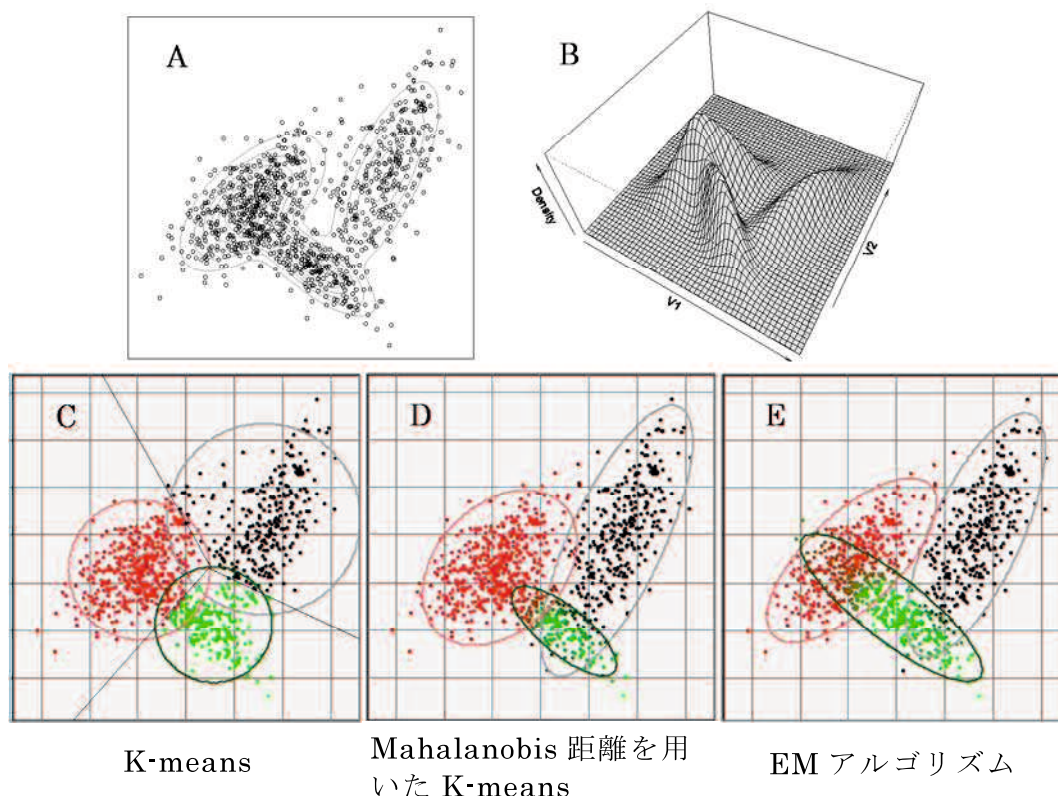


図 2-7 K-means 法と Mahalanobis 距離を用いた K-means 法および EM アルゴリズムのクラスター分析の関係

境界線を見ると、K-means 法ではクラスターが直線によって区分され方法に対して、Mahalanobis 距離を用いた K-means 法 や EM アルゴリズムはクラスターに一致した領域を形成している。特に、EM アルゴリズムはオーバーラップに対して適切な処理が行われている。

EM アルゴリズムは、Expectation と Maximization の 2 つのステップから成り立っている。E (Expectation)ステップ:はじめに適当な初期値を設定し、定義される完全データの対数尤度の条件付き期待値を計算する。この式は、目的である 2.16 式の対数尤度 $LL(\theta)$ の最大化を達成する最尤推定値を導くものである。続いて M ステップでは、E ステップで求めた尤度の期待値を最大化するようなパラメータを求める。

具体的なフローは以下の通りである

- 1) パラメータの初期値を適当に設定する。
 (一次元のときのパラメータは、平均 μ 、分散 σ^2 、
 二次元のときのパラメータは、平均 μ 、分散 σ^2 、相関係数 ρ を使用する)
- 2) 次の E, M ステップを収束するまで交互に繰り返す。
 E ステップ: 今のパラメータでの、各潜在変数 z_i の期待値 $p(z_i | x_i)$ を求める。
 M ステップ: この期待値を重みとして使った対数尤度 $LL(\theta)$ が最大化する方向に
 重心を移動する。このとき新しい平均 μ_k 、分散 σ_k^2 は(2.18) (2.19)式から求める。
- 3) 上記のステップを繰り返し、パラメータの変化量が十分に微小になるまで繰り返す。

$$p(z_i | x_i) = \frac{z_i}{\sum_j^k(z_j)} \quad (2.14)$$

$$f(x_i | \mu_k, \Sigma) = \frac{1}{\sqrt{2\pi}^M} |\Sigma|^{-\frac{1}{2}} \exp\{-(x_i - \mu_k)^T \Sigma^{-1} (x_i - \mu_k) / 2\} = \frac{1}{\sqrt{2\pi}\sigma} \exp\left(-\frac{(x - \mu)^2}{2\sigma^2}\right) \quad (2.15)$$

<一次元正規分布のときの確率密度関数>

$$\text{対数尤度 } LL(\theta) = \sum_{i=1}^n \ln f(x_i | \mu_k, \sigma_k^2) = -\frac{n}{2} \ln(2\pi\sigma_k^2) - \frac{1}{2\sigma_k^2} \sum_{i=1}^n (x_i - \mu_k)^2 \quad (2.16)$$

$$f(x, y) = \frac{1}{2\pi\sigma_x\sigma_y\sqrt{1-\rho^2}} \times e^{-\frac{1}{2(1-\rho^2)}\left\{\frac{(x-\mu_x)^2}{\sigma_x^2} - 2\rho\frac{(x-\mu_x)(y-\mu_y)}{\sigma_x\sigma_y} + \frac{(y-\mu_y)^2}{\sigma_y^2}\right\}} \quad (2.17)$$

<二次元正規分布のときの確率密度関数>

$$\mu_k = \frac{\sum z_{ik} x_i}{\sum z_{ik}} \quad (2.18)$$

$$\sigma_k^2 = \frac{\sum z_{ik} (x_i - \mu_k)^2}{\sum z_{ik}} \quad (2.19)$$

(2.15)と(2.17)式は、 μ_k に対する x_i の尤もらしさを表す。

EM アルゴリズムは複雑に入り組んだクラスターに対して優れた手法である。EM アルゴリズムの利点としては、①データに観測できない隠れ変数(潜在変数)がある場合のパラメータ推定に有用な手法であること、②データ数が多いほど精度が上がるという特徴がある。一般的な条件の下で、 $\theta^{(k)}$ が収束すればその収束値は $LL(\theta | y)$ の局所的最大値であることが証明されている¹⁵⁾。図 2-8 は、一次元の混合分布の例を示す。赤で示す混合分布に対し、緑と青の分布を適切に算出している。

一方欠点は、①初期値の影響を受ける、②クラスター数の設定が事前にされていなければならない、③局所最適解は得られるが大域的な最適解とはならず、高密度なデータに対しては極小領域における局所最適解に陥りやすい、④バックグラウンドノイズの影響を敏感受けてしまう、⑤オーバーラップを示すデータの場合に収束までに時間がかかるなどの欠点が指摘され、計算時間は、Mahalanobis 距離を使った K-means 法よりも長い¹⁵⁾。また、EM アルゴリズムの欠点として分散の大きなクラスターに対して分割しすぎてしまう現象や縮約しきれずに分解してしまう場合、縮約しすぎてクラスターを一つにしてしまう現象もあり、特に極端値の影響により、クラスターを作ることができなくなる致命的な欠陥がある。(大域最適解と局所最適解については付録に記す)。

図 2-9 に初期値であるクラスター数の設定の重要性と EM アルゴリズムのノイズの影響度について示す。本来 2 つのクラスターであるものが、ノイズに対しても 1 つのクラスターとして判定しており、ノイズを含んだデータに対しては、有力なクラスター分析法とは言えない。また、高密度なクラスターにおいては、極小領域での最適解が選択される現象が発生し、4 つのクラスターに過分割された結果となった。

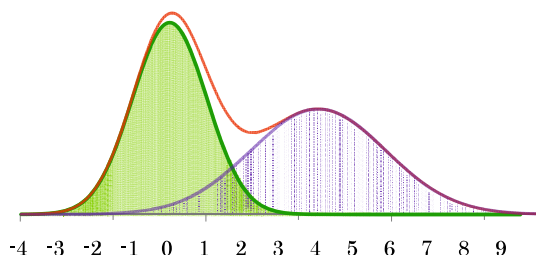


図 2-8 混合正規分布の例

EM アルゴリズムでは、赤線の混合分布から青と緑の分布の算出が可能である。

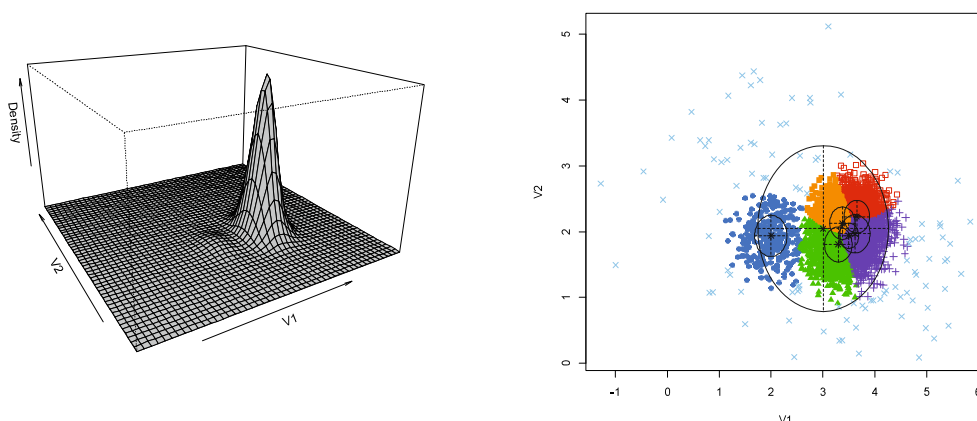


図 2-9 EM アルゴリズムの極小領域での局所最適解とノイズの影響の例

2つの密度と広がりを変えて大きく変えたクラスターに対して、ノイズを発生させて後述の `mclust` を使用して EM アルゴリズムを実行した結果である。EM アルゴリズムでは、ノイズに対して 1 つのクラスターを作り、密度の高いクラスターは分割されて、4 個のクラスターに判定している。

2-4 混合分布モデルにおけるクラスター数の検出

非階層型クラスター分析では、クラスター数は任意に決定するが、ここでは計算による方法を紹介する。分析対象の混合分布モデルを適用する際に、最も重要な問題は、適切に初期クラスターを識別することである。しかし、基本的に分析対象データは、クラスター数が未知であるという問題がある。これを解決する方法として、分散共分散モデルによって推定する方法がある。この方法は、各モデルの条件で、クラスター数を随時変えながら膨大な数の試行錯誤を実施して、ベイズ情報量規準(BIC: Bayesian Information Criterion)を検索し、最も大きな BIC が得られたクラスター数を最適クラスター数としてクラスター分析を行うものである^{17,18)}。したがって、検索には時間を要し、例えばデータ数が 10,000 個を超える場合の計算には、一般的なコンピュータで数分要する。また、検索が失敗した場合は不適切なクラスター形成となる。(BIC については付録に記載する。)

BIC を使ってクラスター数を決定して、EM アルゴリズムを実行するソフトウェアとして R 言語プログラミングで書かれたソフトウェア「mclust」[\[http://www.stat.washington.edu/research/reports/2012/tr597.pdf\]](http://www.stat.washington.edu/research/reports/2012/tr597.pdf)がある¹⁹⁾。mclust は、2次元の混合分布モデルによるクラスター分析を行うために作られたフリーソフトウェアで、BIC 値をグラフで表し、最も大きな値を最適クラスター数とした後、EM アルゴリズムを実行してグラフ化する機能を有している。BIC でのクラスター数の決定は、計算コストがかかるが優れた方法である²⁰⁾。

図 2-10 は、mclust を使用して 3 つのクラスターの立体図を表し、それに対して最適なクラスター数について BIC を求めるための分散モデルを使って検索した結果である。各分散モデルで最大の BIC は異なっている。相関性のあるクラスターに対して EII と VII、EEI といった等分散モデルでは、最大値 4 を示しているが、楕円で異分散、異幅の VVV モデルは適切に 3 クラスターが最大値となっている。このようにクラスター間距離が十分ある場合には、最適なクラスター数を求めることが可能である。経験的には、クラスターの形態が楕円で大きさが異なり、クラスター間距離も異なる VVV 分散モデルが最も EM アルゴリズムでの一致度が高い(表 2-2)。

表 2-2 BIC を求めるための分散モデル

modelNames	分散のモデル	性質
EII	λI	等方, 等分散, 等幅
VII	$\lambda_k I$	等方, 等分散, 異幅
EEI	λA	対角, 等分散, 等幅
VEI	$\lambda_k A$	対角, 等分散, 異幅
EVI	λA_k	対角, 異分散, 等幅
VVI	$\lambda_k A_k$	対角, 異分散, 異幅
EEE	$\lambda D A D^T$	楕円, 等分散, 等幅
EEV	$\lambda D_k A D_k^T$	楕円, 回転同値, 等幅
VEV	$\lambda_k D_k A D_k^T$	楕円, 回転同値, 異幅
VVV	$\lambda_k D_k A_k D_k^T$	楕円, 異分散, 異幅

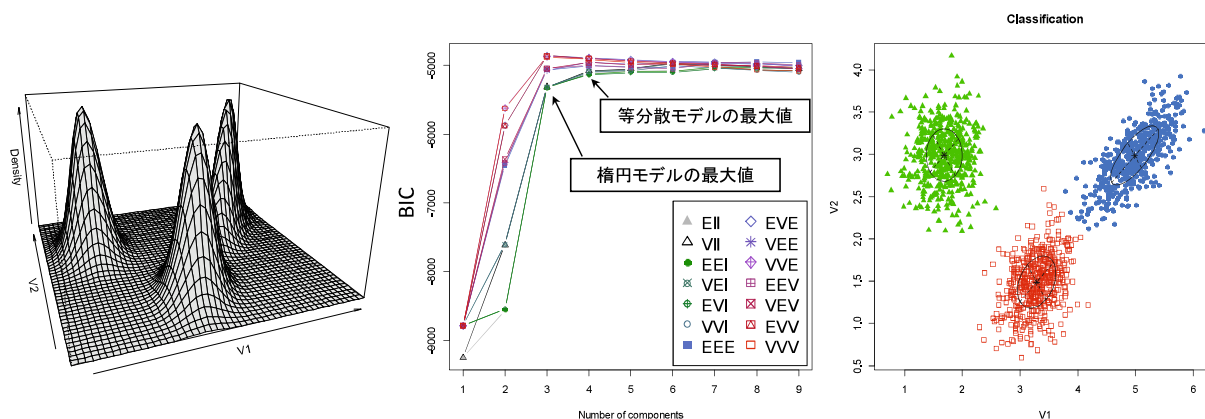


図 2-11 クラスタが明確な場合の BIC によるクラスタ数の検出例
 BIC が分散モデルによって最適クラスタ数が異なっており、等分散モデルの EII, VII, EEI, VEI では最適クラスタ数が 4 であるのに対し、楕円モデルの VVV モデルの最適クラスタ数は 3 であった。楕円で大きさや向きが異なる場合には VVV モデルが最も適している。

2-5 極小領域における局所最適解の影響

クラスタ間に十分な距離がある場合には、クラスタ分析結果は良好な結果が得られるが、EM アルゴリズムの場合には、小さな領域で極端にデータ数が多い場合に、過クラスタを作成してしまうという問題がある。図 2-11 は、データ数がそれぞれ 9000 : 500 : 500 のクラスタである。狭い領域に極端に多くのデータ(9000)を発生させて EM アルゴリズムを実行すると、本来 3 つのクラスタであったものが、データ数 9000 のクラスタは過クラスタ化して 6 つのクラスタ形成をし、合計 8 つのクラスタを作っている。

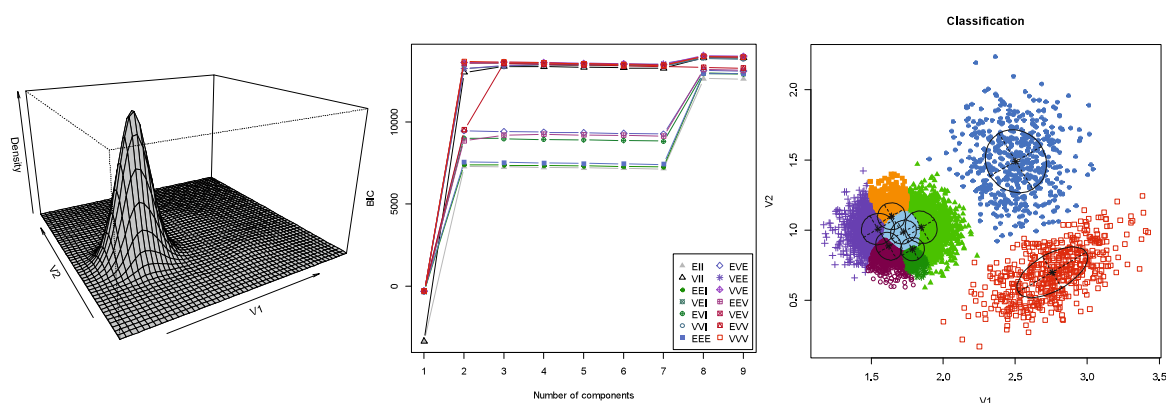


図 2-12 極小領域における局所最適解の影響
 小さな領域で極端なデータ数をもつクラスターに対して、EM アルゴリズムで過クラスターが発生している。

2-6 クラスター分析結果の解釈

一般にクラスター分析は、基本的に探索的なデータ解析手法であり、分割した結果は主観要素や視点要素に基づいて行われていることに注意しておく必要がある。得られた結果から意味を後付けする方法であり、客観的な統計手法ではない。クラスター分析の結果は利用する目的に応じてその妥当性を検証するものである。ただし、後述する白血球分画におけるクラスター分析は、客観的データに基づいて厳格に分類されるべきクラスター分析法と考えられる。

いずれにしても、K-means 法や Fuzzy c-means でユークリッド距離を使ってクラスター分析を行う場合は、クラスター領域の大小や形態とは関係なくクラスター分析が行われることから、クラスター分析の目的としては重心点を検索することが主体であり、クラスターの領域を特定する機能を持っていないと考えられる。したがって、ユークリッド距離を用いたクラスター分析は、臨床領域での使用は限定的であると思われる。

2-7 まとめ

本章では、従来のクラスター分析について解説した。特に大量のデータを取り扱う場合に有利な非階層型クラスター分析について検討データも加えて紹介した。現行のクラスター分析には、以下のような問題点が報告されている。

K-means 法には、最初にランダムに設定される初期値にクラスター分析結果が依存してしまう問題がある⁶⁾。K-means 法の初期値問題は、長年検討され、オリジナルの E. Forgy の方法は、一様にランダムに決められた K 個のクラスターのいずれかに各ポイントを割り当てる手法をとっているが、理論的根拠がないものである。MacQueen らの方法は遠すぎる外れ値を選択

回避するメカニズムはなく、Katsavounidis らの最大距離アルゴリズムは²¹⁾、初期のクラスター重心として最も離れているデータ同士をクラスターの初期値として選択する手法で、クラスター同士の距離が離れたクラスター分析結果を得ることができるが、外れ値の存在により不適切な結果となる⁶⁾。Celebi らも、K-means 法は、クラスター分析結果がランダムに選択される初期値依存し、バックグラウンドノイズや外れ値に非常に敏感に反応してしまうという重大な問題があることを報告している²²⁾。Arthur によって改良された K-means ++法は、最大距離アルゴリズムが最も離れた点同士を初期重心点としたアルゴリズムであるのに対して、互いに近接している 2 つの重心の選択を回避することを目的作られたものにすぎないなどの初期値に関する多くの問題が残っている。

Fuzzy c-means 法は、K-means 法から派生した方法で K-means 法がハードクラスター分析であるのに対して、帰属確率が 0~1 までの間をとるソフトクラスター分析に位置づけられている手法である。通常 Fuzzy c-means がユークリッド距離から算出するのに対して、Mahalanobis 距離を採用した Fuzzy c-means 法は、楕円形(相関関係のある)の 2 次元データには適した方法と思われる。最終的な結果についても初期値問題はあるもののノイズの影響、局所解の問題、極端なデータ数の違いの影響、極端値の影響など多くの改善がなされたものとなっている。

EM アルゴリズムは、理想的な混合分布データを適切に区分する機能を有しているが、他のクラスター分析法と同様に、最初に決めるパラメータ(重心点や分散、相関係数)によって最終的な結果に大きく影響する問題がある。また、未知のクラスター数を既知とし扱わなければならないことも問題である。クラスター数は BIC を使って求める方法もあるが、宮崎らは自由パラメータ数を増やすと混合ガウスモデルの場合にはモデルの当てはまりは良くなり、対数尤度も高くなる。しかし、自由パラメータ数を必要以上に増やすとノイズにもモデルをフィットさせてしまうオーバーフィッティング(過剰適合)が起こると述べている²³⁾。極端に高密度なデータに関して過クラスター化することもあり、さらにノイズや外れ値に非常に敏感で、計算結果が収束しない場合もある²⁴⁾。改良法も考案されているが、未だ解決に至っていないのが現状である。また、今回様々なデータセットについて検討した結果、これまで報告のない密度が高いクラスターが過クラスター化してしまう新たな問題を認めた。臨床分野のデータに対応できるクラスター分析技術が必要である。

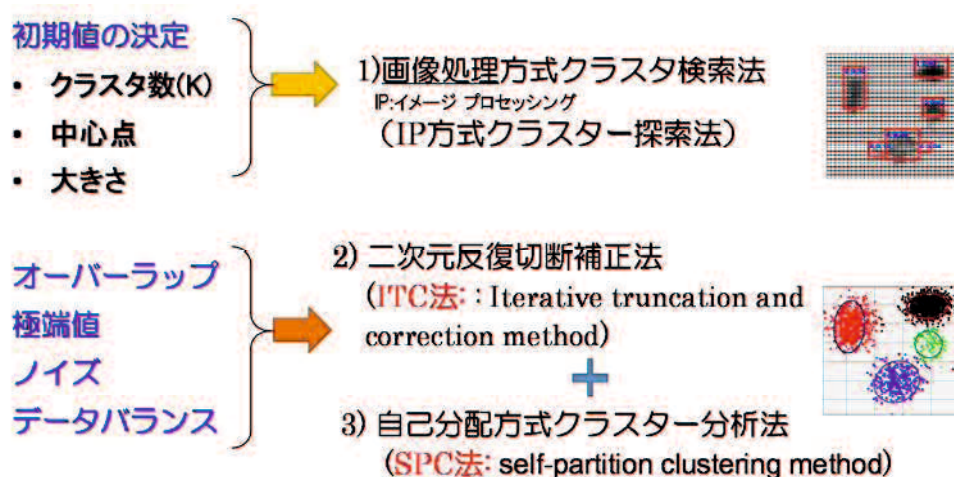
第 3 章

新しいクラスター分析法

- 3-1 イメージ プロセッシング方式クラスター探索法
(Image-Processing (IP) algorithm for cluster search)
- 3-2 自己分配方式クラスター分析法
(SPC: self-partition clustering algorithm)
- 3-3 二次元反復切断補正法
(Iterative truncation-correction (ITC) method)
- 3-4 開発したプログラムについて
- 3-5 比較検証法
- 3-6 まとめ

第 3 章 クラスター解析問題の解決策

現行のクラスター分析技術の問題点として、第 2 章で示したとおり、分類結果が、クラスター初期値に依存すること、極端値やノイズの影響されやすいこと、クラスターの重なり部分の処理が適正に行われないなどがあげられる。実データにおいては極端値やノイズは避けられないものであり、他の問題も的確に打開しなければ臨床検査データへの応用は難しい。今回、新しく開発したクラスター分析法は、このような従来法の問題点を踏まえて、正確な初期クラスターの識別の仕組みとして画像処理方式クラスター探索法を考案した。また、重複データの影響を最小限に抑える二次元反復切断補正法および各データは全クラスターに所属する確率として扱う考え方を導入した自己分配方式クラスター分析法の 3 つの技術を組み合わせた分析手法である。



3-1 画像処理方式クラスター探索法

(Image-Processing (IP) algorithm for cluster search)

これまでのランダム関数を使用して決定していた初期パラメータを的確に推定するために、我々は 2 段階の画像処理アルゴリズム(IP 方式クラスター探索法)を考案した。一般に、画像処理(IP: Image-Processing)は、デジタル化された画像配列情報に、様々なフィルタ関数を当てはめて、いろいろな画像の特徴を捉えたり、元画像を修正したりする技術である。今回この技術を、2 次元配列でマッピングしたものに適用し、データの相対的集積度を求めてクラスターの探索を行った。

現行の非階層クラスター解析では、クラスター数をあらかじめ決めておき、その数に合わせてランダムに初期値の設定が行われるが、その初期値によって分析結果が影響されることをこれまで述べてきた。そこで、表 3-1 に示すフィルタ行列を、2 次元散布図から作成した密度行列 D に適用することで、相対密度(RD: relative density)行列を作成するように設計した。これにより、低密度のクラスターが鮮鋭化される。つぎに、RD 行列の中の値が最大となるブロック(領域)の重心を、トレーサ行列で走査して見つける。見つければ、そのクラスターが占める領域を調べるため、密度が大きく変化しない範囲(80%以上低下しない範囲)で、上下・左右に行列の範囲を拡大する。拡大が終了すると、次に検出されたクラスター領域以外の部分について、密度が大きなブロックを探索し、見つければトレーサ行列でその領域密度を調べて新たなクラスターとする。この一連の操作を反復することで、実在クラスターの数とそれぞれの領域を推定する。IP 方式クラスター探索法の適用で、従来のクラスター分析の大きな課題であった、初期クラスターの数と重心の推定を的確に行うことが可能になる。この処理の具体例として、シミュレーションによる低密度データの変化を図 3-1 に示す。

表 3-1 相対密度作成に使用したフィルタ行列 F

-3	-2	-1	0	-1	-2	-3
-2	0	1	1	1	0	-2
-1	1	2	3	2	1	-1
0	1	3	5	3	1	0
-1	1	2	3	2	1	-1
-2	0	1	1	1	0	-2
-3	-2	-1	0	-1	-2	-3

7×7 の大きさとし、フィルタ値は重心に 5 を与え、周辺に行くに従って円周上に数値を減じるものとした。

IP 方式クラスター探索法の流れ

- 1) データの領域、密度から正方形のブロック数を決定する。
- 2) 与えられた2次元配列 (X, Y) から、密度行列 D を求める。
- 3) 7×7 の画像フィルタ行列(表 3-1)を使って、密度行列 D を RD 行列に置き換える。このとき、3.1 式を使用して密度変換を行う。
- 4) トレーサ行列(スキャン・マトリックス)を使って、RD 行列の最大点を調べ、それを最初のクラスター重心点とする。

- 5) 最高相対密度の位置で、相対密度が 80%以下とまらない範囲で、トレーサ行列を上、左右に広げて、クラスターの領域を調べる。
- 6) 拡大が終了したら、検出されたクラスター領域以外の領域から、RD が最大となる領域を、再びトレーサ行列を用いて探索し、見いだされれば、そのフィルタ行列を確認して新しいクラスターの領域を決定する。
- 7) 6)の操作を繰り返して、初期クラスターの数 K とクラスター領域を決定する。
- 8) 各初期クラスターに対して、その領域に属するデータからクラスターの重心と標準偏差と相関係数を求め、各クラスターの初期パラメータとする。

$$RD[r, c] = \frac{\sum_{u=-3}^3 \sum_{v=-3}^3 D[r-u, c-v] * F[r+3, c+3]}{\sum_{u=-3}^3 \sum_{v=-3}^3 D[r-u, c-v]} \quad (3.1)$$

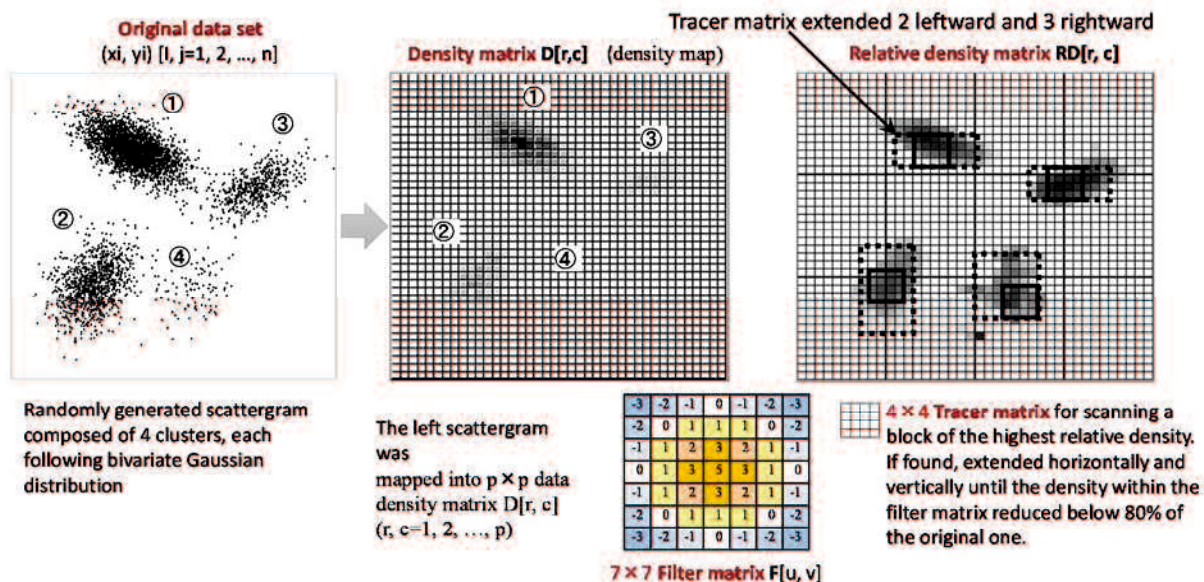


図 3-1 IP 方式クラスター探索法の初期クラスターを検出するための 2 段階の画像処理工程

左側のパネルは、ランダムな 4 つのクラスターを生成した n 個の点 $(X_i, Y_i) [i, j=1, 2, \dots, n]$ からなる二次元散布図である。中央のパネルは、二次元散布図から構築した $P \times P$ の密度行列 D を示す。右側のパネルには、IP 方式クラスター探索法によって得られた RD 行列を表す。低密度クラスターは、コントラストを高めるために 7×7 フィルタ行列 F を使った。フィルタ行列は重心に向かってデータの集積度を走査するために適用した。

右端のパネルの 4×4 の矩形 (実線) は、 RD 行列内の最高相対密度の位置を走査して、最初のクラスター重心を識別したトレーサ行列の位置を示す。この初期矩形は、元の密度の 80% 未満とまらない範囲で、水平方向と垂直方向に拡張される。点線の矩形は拡張後の最終的なクラスターの領域を示す。

3-2 自己分配方式クラスター分析法(SPC: self-partition clustering algorithm)

3-2-1 SPC 法の理論

SPC 法の理論は、次の 2 点である。①各クラスター内のデータ点が二変量ガウス分布に従うものと仮定し、各データ点とクラスター重心との距離には Mahalanobis 距離を採用する。Mahalanobis 距離 D^2 は 3.2 式から算出する。このとき、各データ点のクラスターへの分配率は、 D^2 が自由度 2 の χ^2 分布の上側確率として得ることができる。②各点の、各クラスターへの分配率(S-P: self-partition)値は、 D^2 とクラスターのデータ数で決定される。すなわち、全てのデータが計算に利用されるが、クラスターから遠く離れたデータ点は、 D^2 が非常に小さくなるため、実際の統計計算には含まれなくなる。この点は K-means 法などでは、距離に関係なく全データポイントが計算に利用されることから、考え方が大きく異なる点と言え、データ数を用いる点は従来のクラスター分析とは大きく異なる点であり、独特なものである。

3-2-2 SPC の統計量計算式

SPC 法の計算理論は、各データポイントは、全 K クラスターに所属しており、その第 g クラスターへの分配率は、Mahalanobis 距離 D^2 (3.2 式)とデータ数により決まるという仮定を基礎としている。すなわち、各点 $[X_i, Y_i][i=1, 2, \dots, n]$ の各クラスターとの距離から所属確率を求める。そして、全クラスターについてその所属確率とクラスターのデータ数の積和を分母におくことによって、各クラスター $g[g=1, 2, \dots, K]$ に帰属する確率 $P[g]$ として算出される。さらに、各クラスターのデータ数 $N[g]_i$ と帰属確率 $P[g]_i$ から、i 番目の点のクラスターに属する相対的な量である S-P 値は、3.7 式の $w[g]_i$ を使って算出する。

$$D[g]_i^2 = \frac{zx[g]_i^2 + zy[g]_i^2 - 2 \times r[g] \times zx[g]_i \times zy[g]_i}{1 - r[g]^2} \quad (3.2)$$

$$zx_i = \frac{x_i - Mx[g]}{SDx[g]}, \quad zy_i = \frac{y_i - My[g]}{SDy[g]}; \quad (3.3)$$

$$Mx[g] = \frac{\sum_{i=1}^n w[g]_i \times x_i}{\sum_{i=1}^n w[g]_i} \quad My[g] = \frac{\sum_{i=1}^n w[g]_i \times y_i}{\sum_{i=1}^n w[g]_i} \quad (3.4)$$

$$SDx[g] = \frac{\sum_{i=1}^n w[g]_i (x_i - Mx[g])^2}{\sum_{i=1}^n w[g]_i} \quad SDy[g] = \frac{\sum_{i=1}^n w[g]_i (y_i - My[g])^2}{\sum_{i=1}^n w[g]_i} \quad (3.5)$$

$$r[g] = \frac{\sum_{i=1}^n w[g]_i (x_i - Mx[g])(y_i - My[g])}{\sqrt{(\sum_{i=1}^n w[g]_i (x_i - Mx[g])^2) \times (\sum_{i=1}^n w[g]_i (y_i - My[g])^2)}} \quad (3.6)$$

(M = 平均, SD = 標準偏差, r = 相関係数)

$$w[g]_i = \frac{P[g]_{i \times N[g]}}{\sum_{j=1}^k P[j]_{i \times N[j]}} \quad (3.7)$$

3-3 二次元反復切断補正法 (ITC: Iterative truncation-correction method)

臼井は一次元の反復切断法を考案したが²⁵⁾、市原はその一次元を二次元に発展させ二次元反復切断補正法 (ITC) を考案した²⁶⁾。ITC は、クラスター分析におけるオーバーラップの影響および極端値やノイズの影響を最小限にする方法である。SPC の統計量計算式ではオーバーラップ影響が避けられないため、二次元データ分布の中央部はほぼ普遍的に正規分布であることに着目して、クラスターの中央部 (信頼楕円の内側のデータ) のデータから各クラスターの重心とクラスター内分散を推定するものである。

原理は、極端値やオーバーラップの可能性のある二次元データに対して、S-P 値を利用して中央部の 80% 程度の確率領域部分の基本統計量 (平均値、標準偏差および相関係数) から計算結果が安定するまで繰り返して、本来のクラスターの推定を行うものである。クラスターの周辺データは切り捨ててしまうため、切り捨てた部分を補正係数で調整する。具体的には、80% で切断した場合には、収縮率の逆数に相当する補正係数として、計測された標準偏差 SD に対して 1.29 倍したものを適応する。収縮率の逆数については 3-3-2 補正係数の算出法を示す。

また、すべてのクラスター統計量が更新されたとき、データ数が全データの 1/1000 になった場合は、そのクラスターは最も近いクラスターにマージする仕組みとした。マージ機能によって、IP 方式クラスター探索法で決定された初期クラスター数にとらわれることなく、柔軟なクラスター分析法となっている。この、クラスターを更新するプロセスは、対数尤度和で評価し、データが安定するまで繰り返す。

ITC は、SPC と組み合わせることによって、S-P 値から、個々の点がどのクラスターに近いかを確率的に計算するものである。そして、その確率は各クラスターに所属数を算出することを可能にする。なお、切断する領域は 80% に固定されたものではなく、任意に変更可能な柔軟性をもっている。

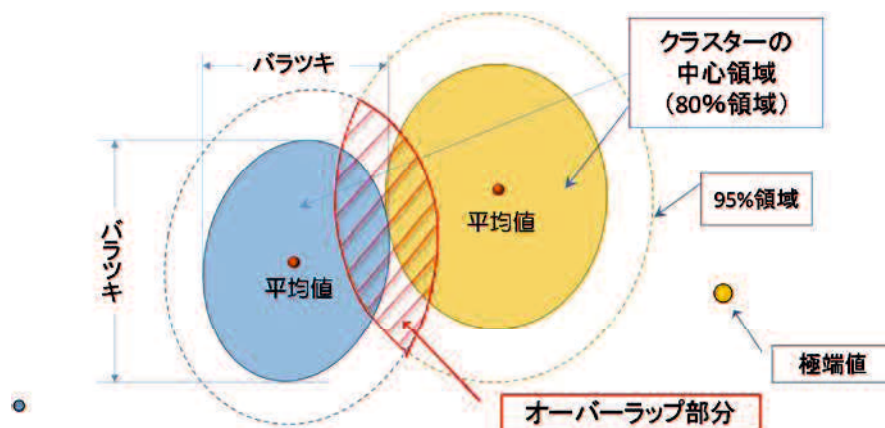


図 3-2 二次元反復切断補正法理論の概念図

点線はデータが存在する 95%信頼区間を示し、濃い色の部分は 80%信頼区間を示す。オーバーラップする領域があっても、濃い色の 80%領域から全体を推定するため、オーバーラップの影響を受けずにクラスター分析が可能である。

3-3-1 SPC/ITC 法の手順

(80%領域で推定する場合)

SPC/ITC 法は、SPC による各データポイントの帰属に関する概念と ITC によるオーバーラップ等の問題を最小限にする方法を合わせたものである。これにより、これまでのクラスター分析の問題点を解決する。SPC/ITC 法とこれまでのクラスター分析法の違いは、K-means では 1 つのデータが必ずいずれかのクラスターに属する事を前提としているのに対し、1 つのデータが全クラスターに帰属し、帰属する確率とデータ数によって各データポイントの S-P 値を算出する方法である。また、EM アルゴリズムが尤度からの確率を採用しているのに対して Mahalanobis 距離 D^2 とクラスターのデータ数から求めるという違いがある。Mahalanobis 距離 D^2 は、臨床検査分野における各種のデータに適用可能な距離であり、完全に一致しないデータ分布の場合でも、変数変換などのデータ処理を行う事で Mahalanobis 距離 D^2 に一致させることが可能である。

SPC/ITC 法の切断率を 0.2 とした時の計算手順例

- 1) IP 方式クラスター探索法によって仮のクラスターである初期パラメータ(重心点、分散、相関係数)の設定を行う。
- 2) 各点からの暫定的クラスター重心点までの Mahalanobis 距離 D^2 の算出。
- 3) D^2 を基に自由度 2 の χ^2 の上側確率から各点のクラスターに対する帰属確率 $P[g_i]$ を算出。
- 4) ノイズの影響を低減するために帰属確率が 0.001 以下のものにはフィルタをかける。

- 5) 二次元反復切断: 帰属確率 $P[g]_i$ が 0.2(20%)以下のデータを切断処理し、残った帰属確率が0.8(80%)以上のデータから仮の重心点、分散、相関係数を再計算する。
- 6) 自己分配値(S-P 値)の算出(帰属確率 $P[g]_i$ とデータ数 $N[g]_i$ から S-P 値を算出)。
- 7) 補正法の実施: 全データの S-P 値を算出して各クラスターの割合を計算する。データ数補正と標準偏差補正を行う。標準偏差補正は、元の全体像を求めるため帰属確率 0.8 で切断した場合は 1.29 倍で補正することで、基データの帰属範囲が算出できる(26, 27)。
- 8) データサイズがある一定数未満になった場合には、そのクラスターはマージされる。
- 9) 2)~8) の計算を、重心が動かなくなるまでまたは対数尤度和(ssLL)が特定の値になるまで繰り返す。

$$ssLL = \sum_{i=1}^n \sum_{g=1}^k \log(N[g] \times P[g]_i) \quad (3.7)$$

3-3-2 補正係数の算出

補正係数は、次のような方法で求めた。

- ① 相関性のある 2 変量(x, y) 正規性乱数を 100 万セット発生させ、x と y の平均と標準偏差を算出する。
- ② 各点の Mahalanobis 距離を計算し、判別値として自由度 2、確率 P の $\chi^2(2, P)$ を使って、判別値よりも大きな値を除外する。
- ③ 除外後のデータから x と y の標準偏差を算出し、その逆数を標準偏差の補正係数とした。また、除外後のデータ数と基のデータ数からデータ数の補正係数を求めた。確率を変えながら作図すると図 3-4 に示すようなグラフになる。

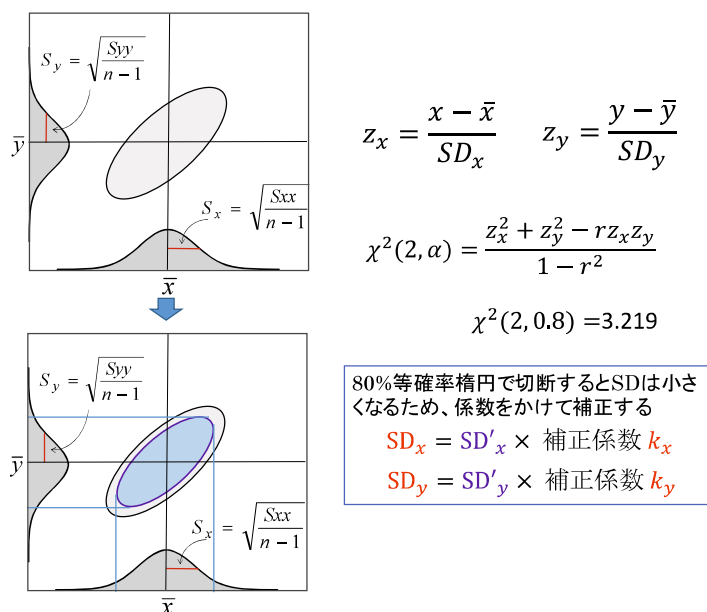


図 3-3 反復切断補正法の算出方法

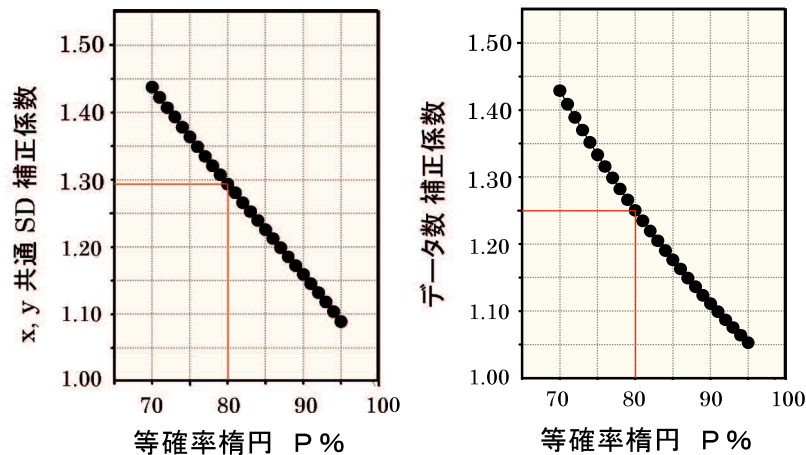


図 3-4 二次元反復切断補正法の補正係数

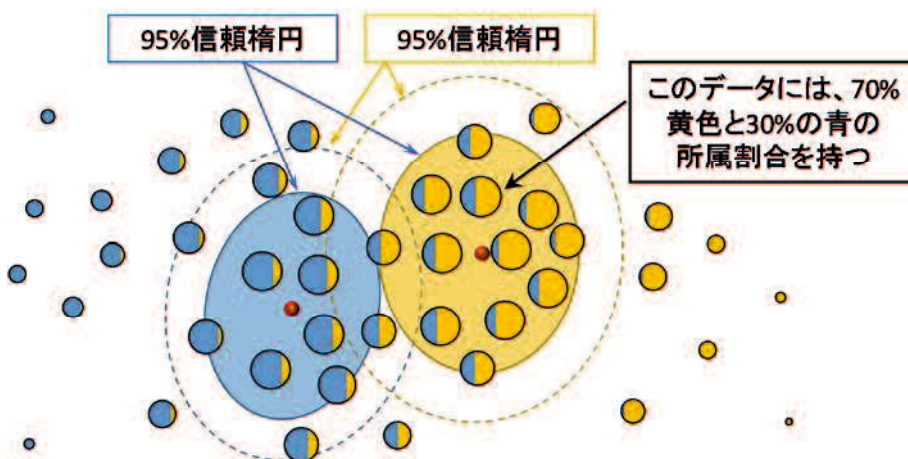


図 3-5 SPC/ITC 法の概念図

楕円は IP 方式クラスター探索法から求めた初期パラメータによって作ったクラスターである。破線は 95%、実線は 80%の信頼楕円を示し、重心とクラスターの広がりから算出した。一方、正円は各データポイント示し、正円の大きさはクラスター重心点からの Mahalanobis 距離 D^2 から割り出した帰属確率 $P[g]_i$ で、分割した色の面積は S-P 値を表している。すなわち、赤で示したクラスター重心点に近いデータほど大きな円となり、離れるにしたがって小さくなる。また、分割された色は、Mahalanobis 距離 D^2 と各クラスターのデータ数から導き出された S-P 値で両クラスターのデータ数が同じ場合、クラスターの間中位置にあるものは半分ずつ所属する事になるが、どちらかのクラスターに偏るとその面積は変わる。

ITC は、初期パラメータから算出された 80%信頼楕円領域のデータから再度重心と標準偏差を計算し、80%に縮小されたデータを補正係数で調整して元に戻すことの繰り返しによって真のクラスターを検索する。

3-4 開発したプログラムについて

プログラム言語 C# (Visual Studio2015. マイクロソフト社)を使用して検証環境作成した。本プログラムは以下のような機能を有する(図 3-6)。

- 1) 指定した数のガウス分布クラスターの自動作成または指定ファイルからのデータ呼び込み
- 2) 作成したクラスターの平均値、標準偏差、相関係数、データ数を表示
- 3) 作成したクラスターの密度 D の描画
- 4) IP 方式クラスター探索法による RD 行列の結果と初期領域の作画
- 5) 各種アルゴリズム(ユークリッド距離, K-means, Mahalanobis 距離, Fuzzy c-means, SPC/ITC 法, EM アルゴリズム)の選択と計数値の設定
- 6) 各種アルゴリズムの計算結果表示(尤度、データ数)、描画表示
各点は、各クラスターに属する S-P 値の大きさによって色分けされ、クラスターの間位置にある場合は、両クラスターの間的な色になるよう調整される。
- 7) クラスターの変化の過程を観察できるように、1ステップずつ処理が進むようにしている。

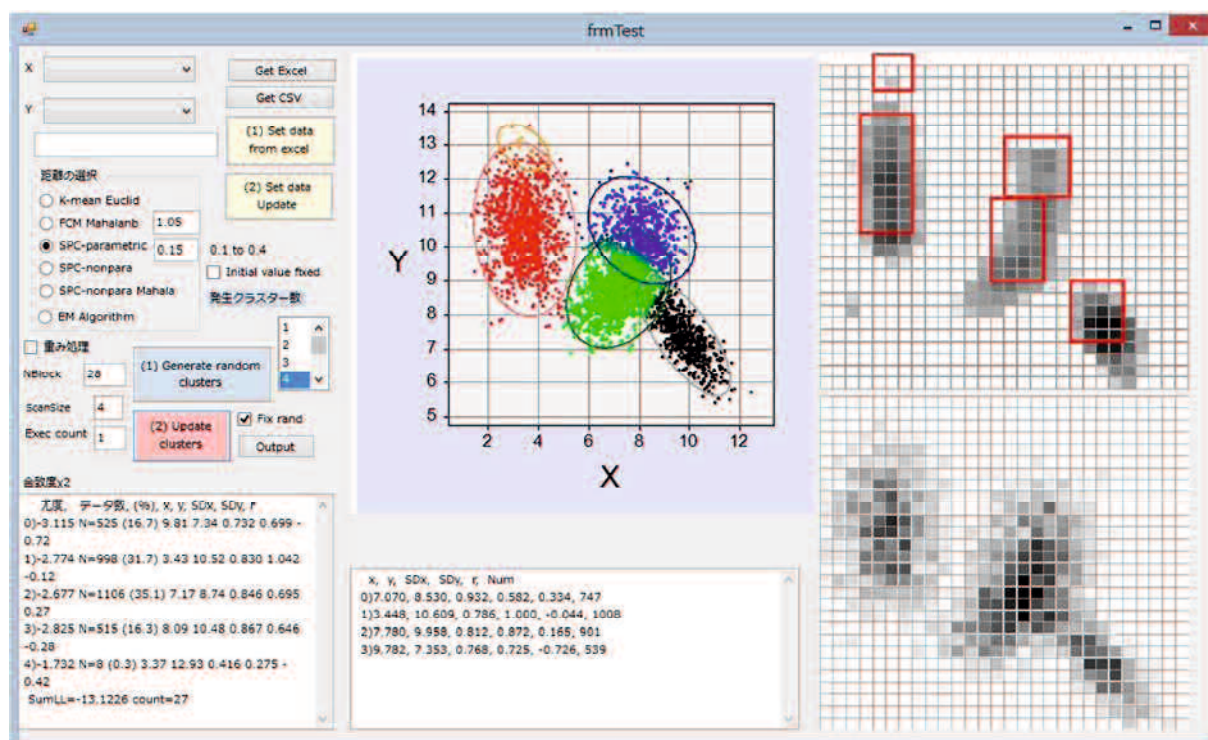


図 3-6 開発したプログラムの実行結果

上段は、自動作成した二次元散布図に対して密度行列、RD 行列を作成し、クラスター分析を実施したところ。

3-5 各クラスター分析法の比較検証法

SPC/ITC 法との比較対象法は、非階層クラスター分析法の基本となる K-means 法とその発展型である Fuzzy c-means 法、混合分布に対して適切な処理を行うとされる EM アルゴリズムの 3 法を中心に評価を行った。特に、EM アルゴリズムは、臨床検査の中に含まれる多くの混合データ処理法として注目されるクラスター分析法である。EM アルゴリズムは、開発したソフトウェアによる方法と R 言語プログラミングの `mclust` の 2 法を使用してシミュレーションデータについて検討を行った。EM アルゴリズムを 2 つの方法で検証する理由は、EM アルゴリズムは初期値の影響を受けて結果が大きく異なるため、開発した IP 方式クラスター探索法を使った初期値の結果と `mclust` の結果を確認するために使用した。

3-5-1 シミュレーションデータによる IP 方式クラスター探索法の性能評価

人工的に作成した正規分布からなる 2 次元クラスターに IP 方式クラスター探索法でのクラスターの検出性能を確認した。クラスターの広がり(クラスター内分散)変えずにデータ数(密度)のみ変化させた場合について図 3-7 に示す。異なるデータサイズの 3 つのクラスターからなるデータセットは、(1)がデータ数 200, 200, 200 で、(2)が 200, 400, 600、(3)が 100, 300, 900 である。元の密度の違いに関係なくマッピングされた RD 行列は、いずれの場合も近似するものとなった。また、赤い四角で囲まれた初期クラスター領域は 3 つのクラスターを的確に捉えていた。識別されたクラスターのパターンは多少異なるが、SPC/ITC 法のマージ機能により最終的には 3 つのクラスターの同定につながった。

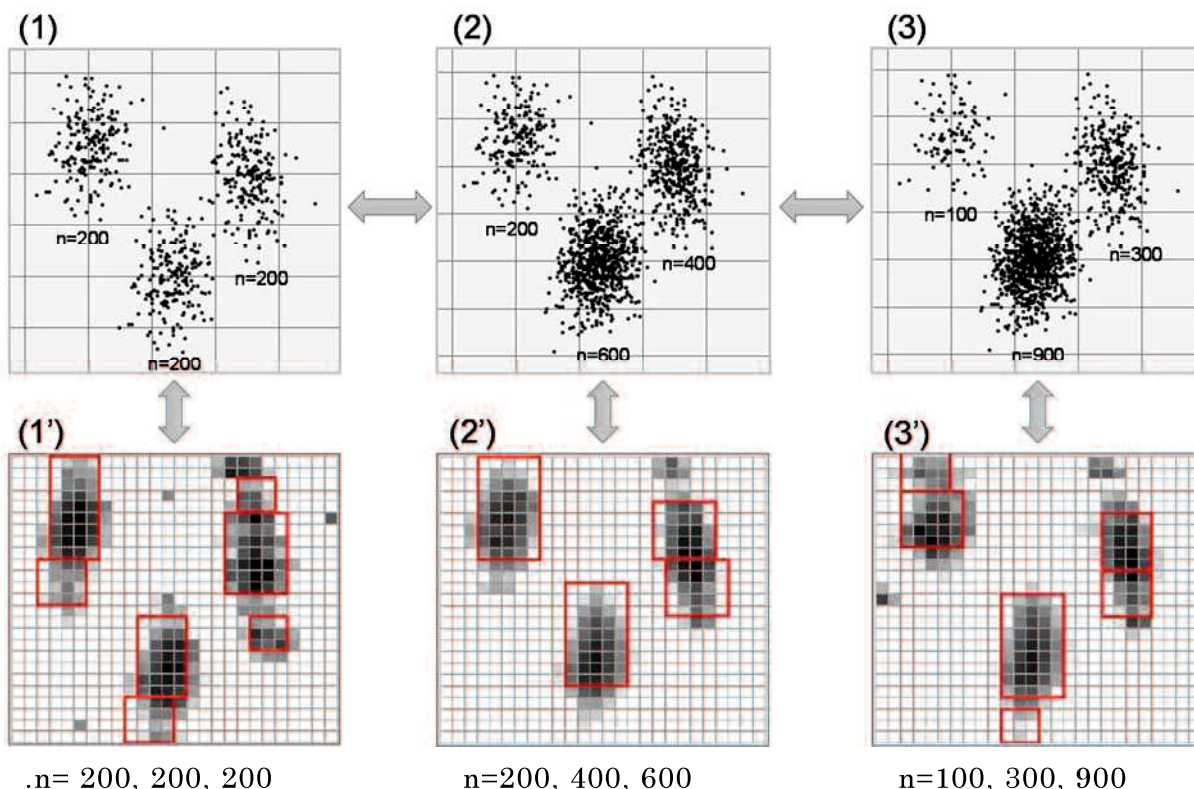


図 3-7 初期クラスターを識別するための IP 方式クラスター探索法の性能
 上段にデータ数（密度）の異なる 3 つのクラスターからなるデータセットを用意し、下段に IP 方式クラスター探索法による RD 行列の結果と赤枠で捉えたクラスター領域を示す。

3-5-2 近接するクラスターにおける SPC/ITC 法の性能評価

図 3-8 は、クラスター間距離とデータ数（密度）を変化させてクラスター分析性能をみたもので、同じ広がりをもつ相関係数 $r = 0.7$ の 2 つのクラスターを示す。クラスターのバラツキは短軸方向で $SD = 1.0$ とし、90%の時の信頼楕円 (CEs: confidence ellipses) では重心点間距離は $3.44 SD$ となり、80%では $2.88 SD$ 、70%では $2.48 SD$ 、60%で $2.17 SD$ となる。データ数が同数 (500:500) の場合、2 つのクラスター距離の 70%信頼区間に換算すると重心間距離として標準偏差の 2.48 倍まで接近させても 2 つのクラスターを判別していた。ただし、60%にすると 2 つのクラスターはマージして1つのクラスターとなった。一方、データ数を 100:1000 のアンバランスにすると 80% ($2.88 SD$) まで判別できたが 70% ($2.48 SD$) では 2 つのクラスターを識別することが出来ずマージしてしまった。

EM アルゴリズムでは、同じデータセットを用いて調べたところ、SPC/ITC 法とは異なり、8 例全てで 2 つのクラスターを識別することに成功した。これは、元来混合分布に対応する EM アルゴリズムは、初期に設定したクラスター数が減ることなくクラスター形成をする方法であるため、2変量ガウス分布を表す純粋なクラスター確実に分離する結果となった(図 3-9)。



図 3-8 SPC/ITC 法のクラスターのオーバーラップの影響

上段に、同じデータサイズ(500 : 500)と相関係数 $r=0.7$ の楕円形状を持つ 2 つのクラスターを示す。クラスターを可視化するために描かれた実線の楕円は、95%信頼楕円 (CEs) で破線の楕円は 80% CEs を示す。下段は、データサイズを 1000 : 100 に変えてクラスターの識別能を見たものである。SPC/ITC 法は、クラスターが近接した場合にマージする機能があるために、極端に近接した場合にはクラスターが一つになった。

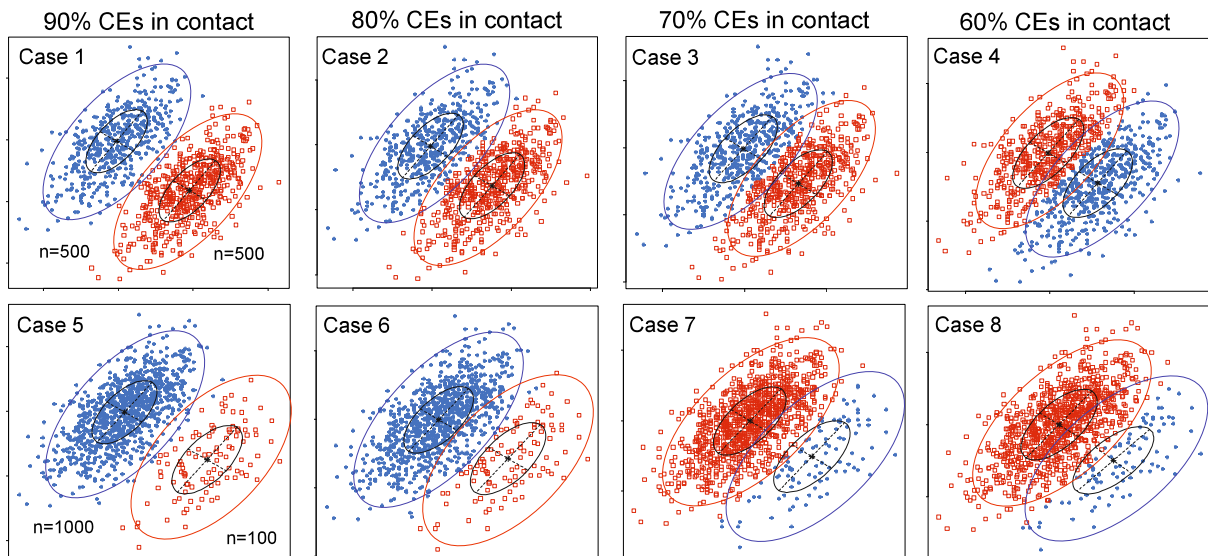


図 3-9 EM アルゴリズムによるクラスターのオーバーラップの影響

図 3-8 と同一データについて `mclust` を使って EM アルゴリズム実施した結果を示す。EM アルゴリズムは元来マージする機能を持たないことから、SPC/ITC のように接近したクラスターであっても 2 つに分離していた。

3-5-3 近接するデータ数(密度)や広がりを変化させたクラスターに対する各種クラスター分析法の比較

クラスターのデータ数やクラスター内分散を大きく変化させた場合の各種クラスター分析法について比較した(図 3-10)。各パネル内の左上の近接する 2 つのクラスター①は、同じデータ数で、クラスター内分散も同じくし、クラスターの関係性の向きも同方向となるようにした。パネル内の右上の 2 つのクラスター②は、データ数は同じであるが、x 軸 y 軸方向ともにクラスター内分散を 2 倍にしたものである。また、左下のクラスター③は各クラスターの向き、分散を変え、データ数の比も 7 倍異なるものにした。元データ(パネル B)の密度や分散をわかりやすくするためにパネル A に元データの立体図を示す。パネル C は密度、パネル D は RD 行列と初期設定値を表す。クラスターが適切に検出できるように階級数を調整した。

各クラスター分析法の結果を図の右側に示す。パネル E は K-means で、ユークリッド距離からクラスターに属するか否かを決定しているために各クラスターの境界面は直線となっている。注目するところは、左下にあるクラスター数が本来 2 つの設定であるが、初期値に 3 クラスターとなってしまったために、3 つのクラスターに分けられていることである。初期値の設定が最終結果まで影響している。Mahalanobis 距離を使った Fuzzy c-means(パネル F)は、クラスターを楕円として識別しているが、K-means と同様左下にあるクラスターを初期値の影響を受けて 3 クラスターとなっている。また、EM アルゴリズム(パネル G)では近接していても比較的低密度のクラスターの識別能は良好で 3 つの混合分布として捉えている。ただし、同じデータについて R 言語プログラミングの mclust で分析すると高密度データで過分割された結果となった(図 3-11)。これは、EM アルゴリズムが極小領域で局所最適解に陥りやすいために起こった現象である。パネル H は SPC/ITC 法で、左下のクラスター③が最終的にマージして 2 つのクラスターとしての的確に捉えている。

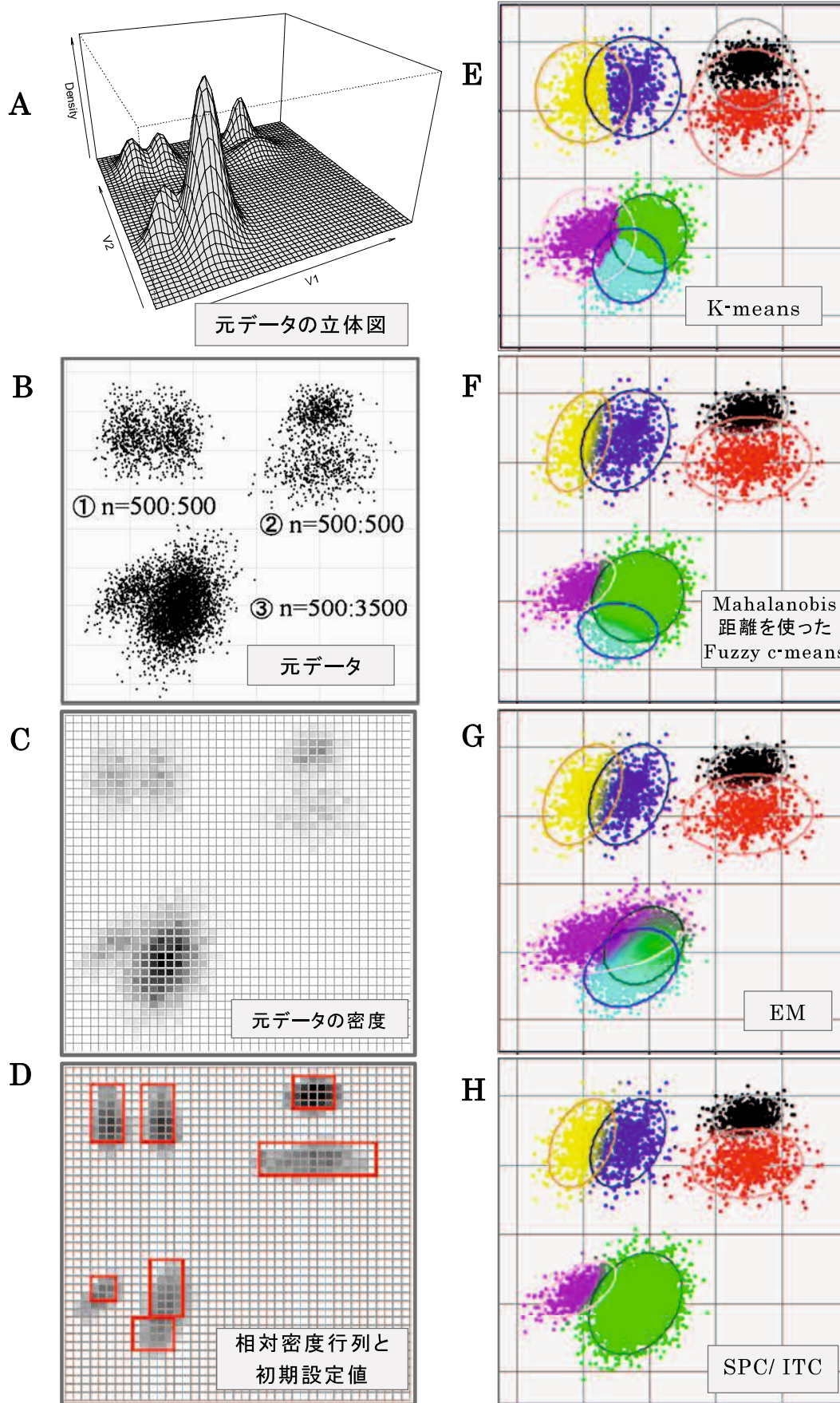


図 3-10 近接するクラスターにおける各種クラスター分析法の比較

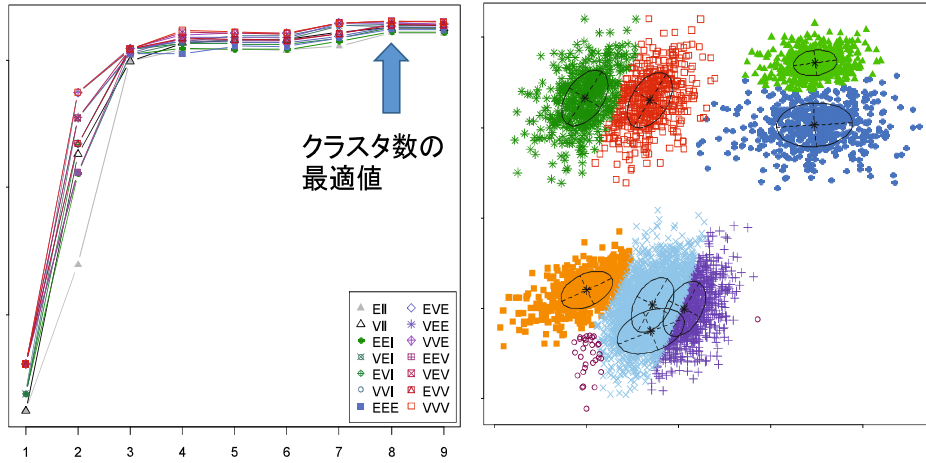


図 3-11 R 言語プログラミングの mclust で分析結果

左のパネルは、最適クラスタ数の計算結果を示し、右のパネルは図 3-10 と同じデータについて、R 言語プログラミングの mclust で分析した結果を示す。高密度クラスターが過分割されている。

3-5-4 EM アルゴリズムと SPC/ITC 法のバックグラウンドノイズの影響

血球自動分析装置における白血球分画などの実データでは、死細胞や凝集、残渣の影響によってバックグラウンドノイズが多く含まれ、クラスター分析では、この影響をなくすことが課題となる。そこで、バックグラウンドノイズを含んだシミュレーションデータでの影響について検証を行った(図 3-12)。データサイズは、それぞれ $n=500$ の 2 つのクラスターに、一様分布のバックグラウンドノイズを 5%から 30%まで変化させて EM アルゴリズムと SPC/ITC 法の違いを見た。

EM アルゴリズムでは 5%ノイズの時点からノイズに敏感に反応して、ノイズ成分をクラスターとして識別していたのに対して、SPC/ITC 法では 20%ノイズの状態まで的確にクラスターを捉えていた。ノイズの影響を受けない技術として ITC 法が大きく貢献していると思われ、このことはバックグラウンドノイズが多く含まれる実データにおいて SPC/ITC 法の有意性を示すものである。

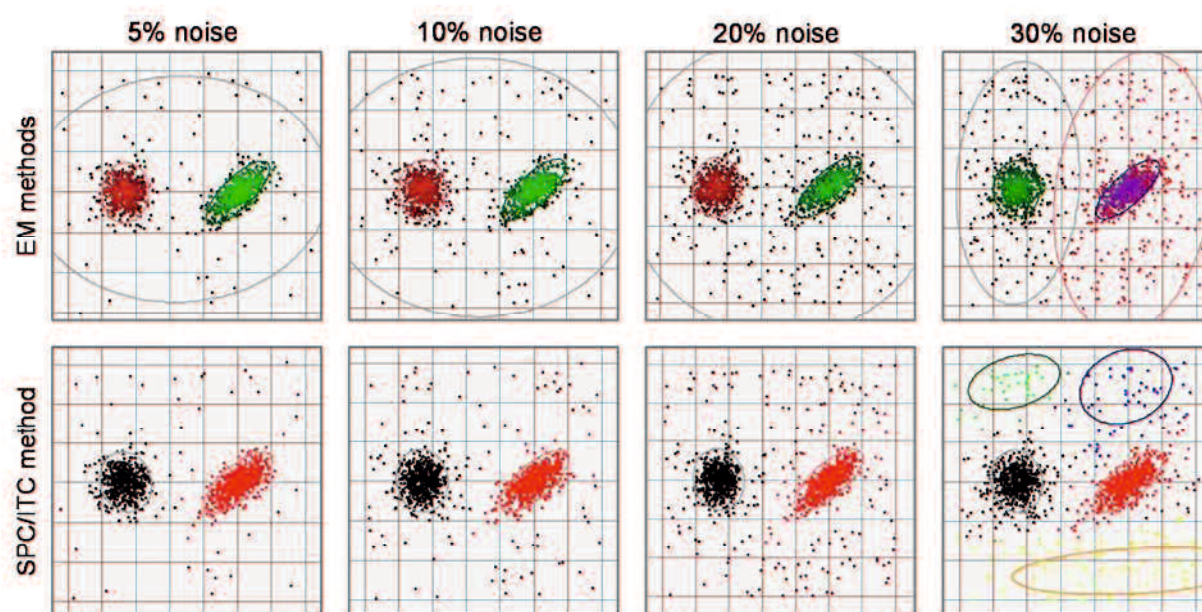


図 3-12 EM アルゴリズムと SPC/ITC 法のバックグラウンドノイズの影響

重複のない $n=500:500$ の相関性のある 2 つのデータセットに、一様分布のバックグラウンドノイズ成分を 5% ($n=50$)、10% ($n=100$)、20% ($n=200$) および 30% ($n=300$) を加えて、EM アルゴリズム(上段)と SPC/ITC 法(下段)で比較した。EM アルゴリズムでは、5~30%の全てにおいてノイズに反応して不適切なクラスターを作成しているのに対して、SPC/ITC 法ではノイズ量 20%まで適切なクラスター分析が可能であった。

3-5-5 白血球分画疑似モデルによる一致率

白血球分画疑似モデルを作成し、各種クラスター分析法での挙動の確認を行った。疑似モデルで検証する目的は、第 4 章の実モデルデータでは、分析機器から出力されたデータが、正確な分類情報を提供していない可能性があるため、ここでのシミュレーション結果が実際上のアルゴリズムの正確性評価となる。

白血球分画疑似モデルの作成に当たっては、血球自動分析装置から出力される健常人モデルを参考にし、モデルデータとして分析装置は堀場製作所の Pentra MS CPR(京都)とシスメックス株式会社の XN-1000(神戸)の 2 機種を参考にした(図 3-13)。各細胞分画は正規分布を想定して疑似データ調整した(図 3-14)。血球自動分析装置の細胞分類に関しては 4-1-4 の白血球分類・フローサイトメトリーの原理を参照。



Pentra MSCPR



XN-1000

図 3-13 今回の検討に使用した血球自動分析装置

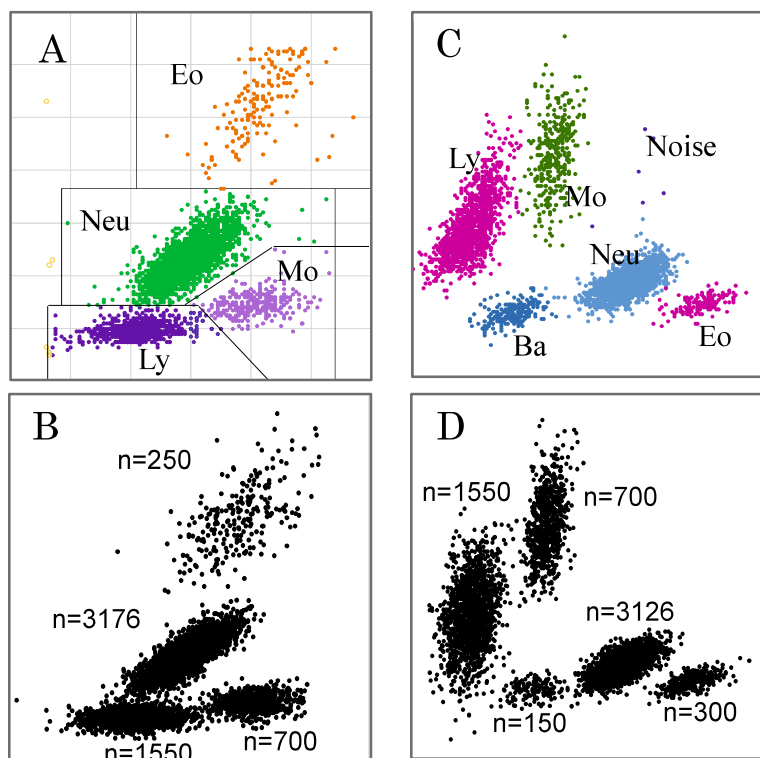


図 3-14 健常人データを使用して分析装置から出力された二次元散布図と人工的に作成した白血球分画疑似モデル

パネル A は堀場製作所の Pentra MS CPR の二次元散布図、パネル C はシスメックス株式会社の XN-1000 の二次元散布図を示す。Pentra MS CPR では、好中球(Neu)、リンパ球(Ly)、単球(Mo)、好酸球(Eo)の 4 クラスターが二次元散布図上に表示されている。XN-1000 では Pentra MS CPR の 4 クラスターに加えて、好塩基球(Ba)領域を加えた 5 クラスターが二次元散布図上に表示されている。パネル B と D は、パネル A, C を参考に白血球分画疑似モデルを作成した結果である。

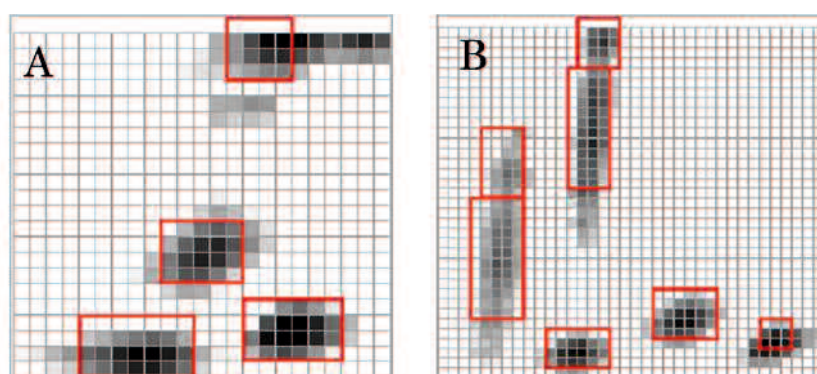


図 3-15 白血球分画疑似モデルの RD 行列と IP 方式クラスター探索法によるクラスター領域

パネル A は Pentra MS CPR の疑似モデルに IP 方式クラスター探索法に適応して 4 つのクラスター領域が検出されていることを示している。パネル B は XN-1000 の場合を示し、本来 5 つのクラスターであるが、7 つのクラスターが検出されている。

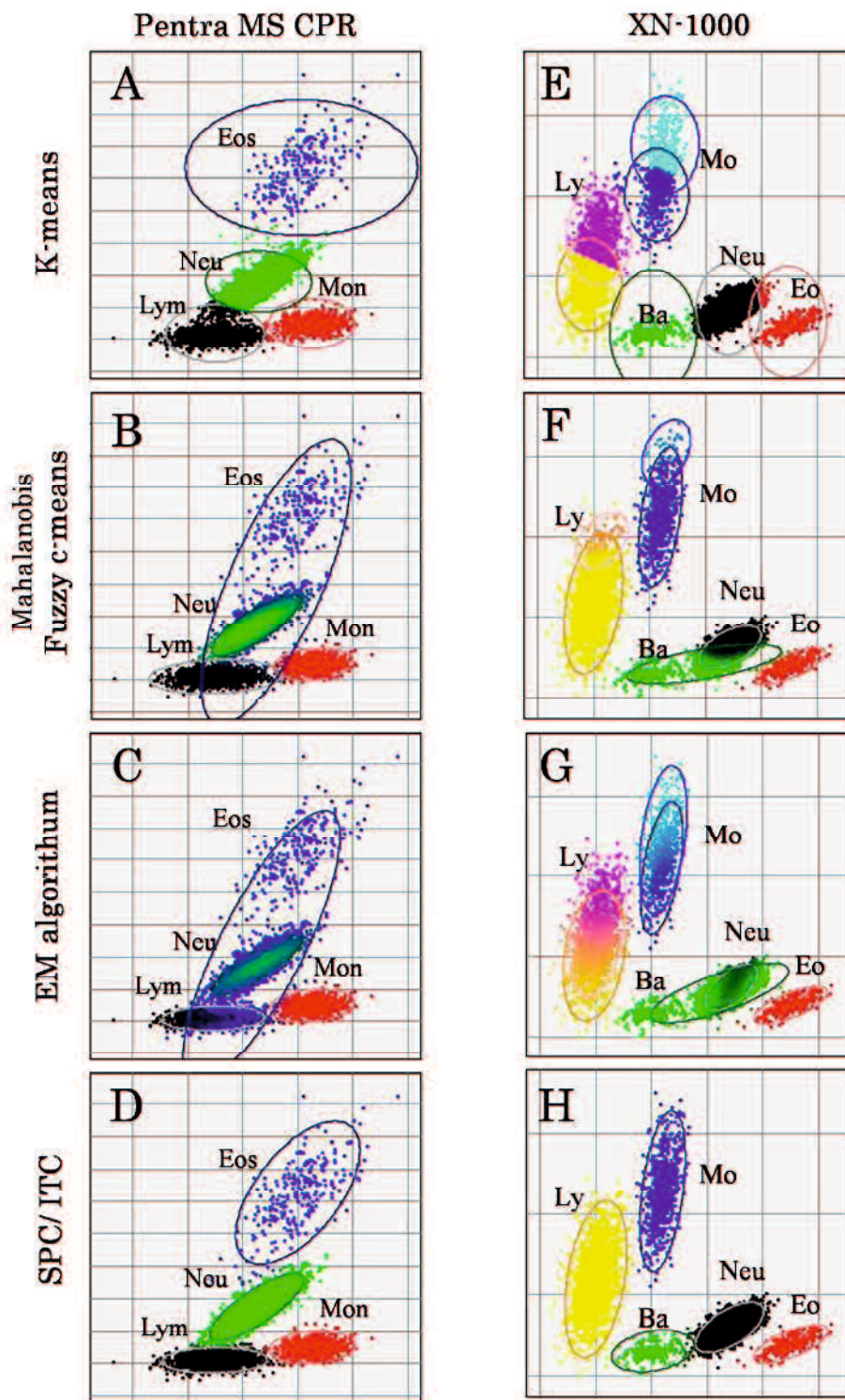


図 3-16 白血球分画疑似モデルにおける各種クラスター分析法の検証結果

表 3-2 白血球分画疑似モデルによる各クラスター分析法の誤差率

Pentra MSCPR	Neu	Ly	Mo	EO
設定値	3176	1550	700	250
k-means	3006 -5.4%	1724 11%	714 2.0%	234 -6.4%
マハラノビス距離	2796	1522	683	676
Fuzzy c-means	-12%	-1.8%	-2.4%	170%
EMアルゴリズム	2400 -24%	1416 -8.6%	703 0.4%	1157 363%
SPC/ ITC	3144 -1.0%	1535 -1.0%	712 1.7%	253 1.2%

下欄は設定値に対する誤差%を示す

XN-1000	Neu	Ly	Mo	Eo	Ba
設定値	3126	1550	700	300	150
k-means	3077 -1.6%	1501 -3%	711 1.6%	332 11%	165 10%
マハラノビス距離	2208	1534	688	289	1070
Fuzzy c-means	-29%	-1.0%	-1.7%	-3.7%	613%
EMアルゴリズム	1851 -41%	1556 0.4%	702 0.3%	298 -0.7%	298 99%
SPC/ ITC	3118 -0.3%	1554 0.3%	698 -0.3%	299 -0.3%	154 2.7%

下欄は設定値に対する誤差%を示す

背景が黄色い部分は、同系のクラスターが2分されている細胞群の値を合計した数値である。

IP方式クラスター探索法で健常人の白血球分画を想定した白血球分画疑似モデルによるクラスター分析結果では、初期値の段階で Pentra MS CPR モデル はシミュレーションで作成した4つのクラスターが適切に検出されているのに対して、XN-1000 モデルでは長径方向に長いクラスターが2つに分けられた。本来シミュレーション上では5つのクラスターであるべきだが、長径を2分する形になったために、合計7つのクラスターとなった(図 3-15)。

図 3-16 は白血球分画疑似モデルの各クラスター分析法の結果を示す。ユークリッド距離に基づく K-means では、Pentra MS CPR モデル、XN-1000 モデル共にシミュレーションで作成したクラスター領域の形態に適合した形に検出できていない。Mahalanobis 距離を使った Fuzzy c-means は、Pentra MS CPR モデルにおいて好酸球(Eo)領域が拡大して好中球(Neu)領域やリンパ球領域にまで及んでいる。XN-1000 モデルでは好塩基球(Ba)領域にあるはずのクラスター領域が大きく好中球(Neu)領域にずれていた。このように Fuzzy c-means は、データ数が少なくクラスター内分散が大きいクラスターが他のクラスター領域を囲い込んでしまう傾向があった。EM アルゴリズムも Fuzzy c-means と同じ傾向が認められたが、オーバ

まとめ

ーラップした分布においても初期クラスター数を維持したクラスターを作成するため、2重円のような形になっている。これに対して、SPC/ITC法は近接するクラスターを融合する機能があるため、作成した白血球分画疑似モデルに完全にフィットし、余分なクラスター形成なく適切な処理がされていた。

誤差率表 3-2 から見ると、Pentra MS CPR モデルでは、リンパ球(Ly)と単球(Mo)は、どのクラスター分析法でも誤差の少ない結果であった。これは他のクラスターの影響が少ない所に位置していることと密度もあるためと思われる。好中球(Neu)については、Fuzzy c-means と EM アルゴリズムにおいてバラツキが大きく密度が小さい好酸球(Eo)の影響を受けて、一部取り込まれてしまうためにマイナス誤差が発生した。これに対して SPC/ITC 法は最大誤差の Mo であっても 1.7% であり良好な結果を示した。XN-1000 モデルでは、Fuzzy c-means や EM アルゴリズムで好塩基球(Ba)が Pentra MS CPR の好酸球(Eo)のように好中球(Neu)の領域まで取り込んでしまい、好中球(Neu)に大きなマイナス誤差が発生している。逆に、好塩基球(Ba)は非常に大きなプラス誤差が発生した。このように現行のクラスター分析法が大きな誤差を発生しているのに対して、SPC/ITC 法は最も大きな誤差%でも好塩基球の 2.7% であり、良好な成績であった。

3-6 まとめ

本章では、新しいアルゴリズム理論である IP 方式クラスター探索法と SPC/ITC 法を紹介し、現行のクラスター分析法と新しいアルゴリズムに基づくクラスター分析法の比較検証を行った。

IP 方式クラスター探索法は、クラスター分析の初期値問題に対する的確な重心とクラスターの分散を検索する仕組みとして、密度によって検索する方法である。この方法は、二次元散布図をブロックのサイズを適切に調整して分割し、分けられた領域内のデータ数を密度と捉え、高密度になったブロックがクラスターの重心であり、クラスターの大きさは高密度なブロックからの変化が少ない領域とする方法である。このとき、画像フィルタリング技術を応用して密度に対する鮮明度を RD 行列に変換して実施した。これまでのクラスター分析のようにランダム関数を使用して仮重心点を設定するクラスター分析に比べ、不適切な初期値が選択されて最終結果に影響する問題が発生しない点で優れた方法といえる。

なお、臨床検査の中での初期値に関しては、白血球分類・フローサイトメトリーやタンパク分画などの場合には、過去の膨大な事例から初期値はある程度推定できるため、教師なし(過去の経験を利用しない)の初期値検索ではなく、ある程度の範囲に初期値を置く半教師ありの方法も考えられる。

クラスター数の問題に関しても SPC/ITC 法では、分析の初期にクラスター数が決定されている必要はなく、IP 方式クラスター探索法によって検出された多めのクラスターから出発して、近接するクラスターがある場合、マージ機能によって近接するクラスターが統合され、クラ

スター数を調整する方法をとっている。この方法によって、多くのクラスター分析法で、不適切なクラスター数が原因での確なクラスター分析となるのを防いでいる。

クラスターへの帰属に関しては、非階層型クラスター分析の代表的存在である **K-means** 法では、個々の点(データ)はいずれかのクラスターに属するハードクラスタリングを前提とした処理法である。比較的新しい考えの **EM** アルゴリズムは、尤度からの個々のデータに対する所属確率を採用してクラスター分析を行うものである。**SPC/ITC** 法ではこれまでのクラスター分析が距離のみで分析を行っていたことに対して、個々の点は各クラスターからの **Mahalanobis** 距離とクラスターのデータ数によって自己分配する概念を導入した。データ数を計算に入れることにより、より分布を反映した分析と考えている。基本の分布として正規分布を仮定した方法であるが、多くの分布は正規分布であり、今回実例として使用した白血球分類データについても正規分布を仮定した方法で十分その機能を果たしている。

また、ノイズについては **K-means** 法や **EM** アルゴリズムなどの現行のクラスター分析では、バックグラウンドノイズなどの僅かなデータによって敏感に影響を受けてクラスターを作ってしまうのに対して、**SPC/ITC** 法では少数のノイズには全く影響を受けずにクラスター分析が行え **EM** アルゴリズムより遙かに優れた結果であった。オーバーラップするクラスターに対しては、**SPC/ITC** は **EM** アルゴリズムほど強靱ではないが、**EM** アルゴリズムにおいて高密度データに対し過分割してしまう現象をマージ機能によって適切な処理が可能であった。

これらのことは、市原の考案したクラスター重心部の情報は、上手くクラスター全体を表しているという考え方から、二次元反復切断補正法を取り入れた効果である²⁶⁾。実際には、そのクラスターが持っている情報の重心部分 70%から 90%程度を使用して全体像を見いだす方法を採用することにより、ノイズの影響や外れ値などの極端値やクラスター周辺のオーバーラップ部分の影響を受けることなく、クラスターが持つ本来の特性を見いだすことに成功している。

計算コストに関しては、**SPC/ITC** 法が **Mahalanobis** 距離から算出することにより、**K-means** 法に比べ大きい **EM** アルゴリズムでの尤度を使用した場合に比べて大幅な削減がなされている。五利江は、混合分布のパラメータの推定において、分布の状況によって最尤法と最小二乗法では計算コスト差が認められることを報告しており²⁸⁾、**SPC/ITC** 法は中間に位置するものと思われる。実験結果よりどのようなデータについて解析を行うかによって計算コストは大幅に異なること明らかである。**Redner** らは、「**EM** アルゴリズムは信頼性の高い大域的収束性があるが、収束は極めて遅い」と述べていることと一致する²⁹⁾。

白血球分画疑似モデルの検討では、データ数の大小やクラスター領域の大小の幅が大きい、ノイズ成分などに対して、**SPC/ITC** 法は正確な分画データが得られた。これまでのクラスター分析において実現できなかった分野に対して、有効な方法であることを確認した。

まとめ

第 4 章

臨床検査分野への応用

- 4-1 泳動分析
- 4-2 自動血球分析装置
- 4-3 実データによる検証
- 4-4 白血球分画の実データによる
各種クラスター解析のまとめ

第 4 章 臨床検査分野への応用

臨床検査分野では、クラスター分析技術を必要とする分野は、タンパク電気泳動などの泳動分析や自動血球分析装置による血球計測、白血球分画の領域などである。しかし、これら分野でのクラスター分析の実績はなく、報告もない。クラスター分析が困難な要因として、正常から病的なパターンまでの様々な病態が存在するために、正確な分類を行うことが困難であることが上げられる。また、白血球分画では、リンパ球と単球、好中球と単球などの細胞構造的に近いものは分類が難しいことがある。ここでは、一次元データである各種泳動検査の原理と問題点および二次元データである白血球分画に用いられるフローサイトメトリー法の原理と問題点を示し、新しいクラスター分析法を用いた白血球分画への応用技術を述べる。

4-1 泳動分析

4-1-1 各種泳動検査の原理

電気泳動の原理

電気泳動は、荷電粒子が電界中を移動する現象で、タンパク質や核酸の主たる分離・分析法となっている。荷電粒子とはペプチド・タンパク質・DNA・RNA などの核酸のことで、タンパク質はアルカリ溶液中で電場を与えると負に荷電して陽極へ移動する性質を利用している。水溶液中では試料が拡散してしまうため、支持体として網目構造を持ったセルロースアセテート膜やアガロースゲルやポリアクリルアミドゲルを用いる。移動速度は、荷電の状況や粒子大きさによって移動速度が異なることを利用している³⁰⁾。この原理を利用して、血清中に存在する約 100 種類のタンパク質を大きく 5 つの分画に分離している(図 4-1, 2)。この分離作業によって分子量決定をはじめ、等電点や純度決定、各成分の定量、精製等に用いられる。

図 4-1 から分かるように、各分画は明瞭に分画されておらず、各分画が連続した分画となる混合分布を示している。

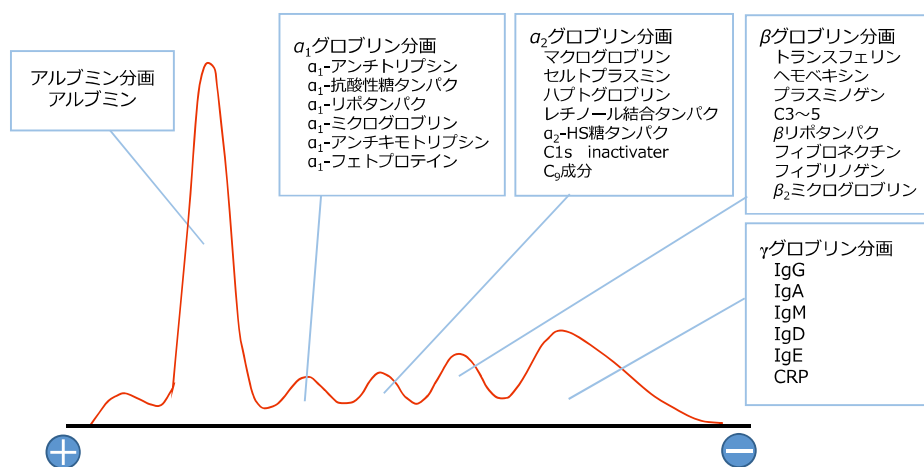


図 4-1 血清タンパク分画のデンストグラムと各分画に属する主要血清タンパク成分
各分画は明瞭に分画されていないことから、混合分布と考えられる



図 4-2 血清タンパク分画像

4-1-2 泳動分析の問題点

タンパク泳動分析は、その境界域を谷値で決定しているが、オーバーラップする成分がある場合は、正確に成分分配がされていないことが予測される。

例としてタンパク電気泳動のような一次元 2 クラスター（緑と青）があるとする（図 4-3）。クラスター（緑）の中心を 10、標準偏差を 3 とし、クラスター（青）は、重心を 20、標準偏差を 5 とした場合、赤で示した混合分布の谷値 16 は混合分布を分配しているが、その成分分配は緑成分が 30% で青成分 70% 含有する。緑と青のクラスターの交点は 14 であり、谷値とは一致していない。K-means 法のように重心からの距離を基準にクラスター分析する方法では、クラスターの中間の 15 は、谷値を基準とした場合のクラスター（緑）に属すると考えられる。しかし、所属率から考えるとクラスター（青）に属すると考えるべきである。したがって、電気泳動のような所属量を検出する場合、ユークリッド距離を使った K-means 法のように、単純に距離によって成分分離することに問題があることが想像される。つまり、境界が明瞭な成分は適切な分離がされていると考えられるが、オーバーラップする場合には両成分は混合データであるとしてデータを処理するべきである。

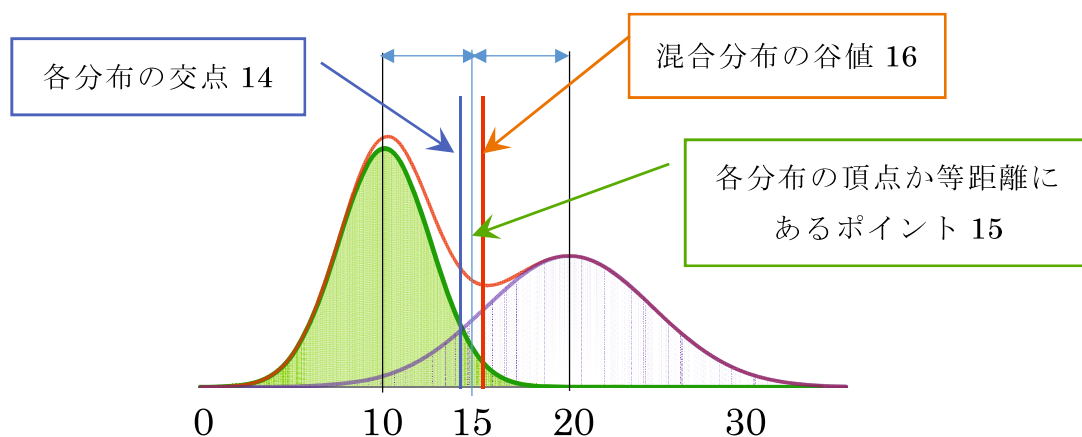


図 4-3 1次元混合分布における境界値の問題点

2 つの一次元クラスターとその混合分布を示す。赤い垂直線は混合分布の谷値であり、タンパク分画に当てはめると、この点を各タンパク成分の境目として処理している。しかし、各分布の交点は、混合分布の谷値とは異なる場所に存在する。

4-1-3 自動血球分析装置の臨床的有用性

臨床検査室における基本的検査として血液学的検査が一般的に実施されている。本検査は、生体の有用な情報として赤血球数、白血球数、血小板数、ヘモグロビン濃度、ヘマトクリットと言った項目ばかりではなく、フローサイトメリー法を使って白血球分画までも小型の機械で行えるようになってきている。白血球分画は、白血球の大きさや染色性、白血球内部の構造情報から細胞分類および白血球の幼若性や異常性までの検出が行われるまでに至っている。この情報を基に、感染症や炎症性疾患の診断、さらには白血病をはじめとする造血器腫瘍の検出と治療経過観察を可能なものとしている。白血病などの診断には、最終的に顕微鏡による血球形態の観察が不可欠なものであるが、血球の数的異常や異常細胞の存在を検出することは、自動血球分析装置に任されている。したがって、的確に細胞分類が出来ることが要求される。

末梢血液中の白血球は形態的に好中球、好酸球、好塩基球、単球、リンパ球の五種類に分類され、各細胞の増加減少パターンによって疾患の推測が可能となる。例えば、好中球が増加している場合は細菌または真菌の感染を疑い、好酸球であればアレルギー性の疾患または寄生虫感染症を考え、単球であれば結核や亜急性心内膜炎などの慢性炎症性疾患、リンパ球が増加していればウイルス感染などを推測する。また、細胞が増加する際には、その細胞の幼弱細胞や炎症反応性細胞なども増加するために、健常な状態では認められない細胞も増加してくる。このように様々な細胞が出現することは、自動分類の困難さ増す要因である。

現在、自動細胞分類法としては、パターン認識の技術を使って細胞の微細構造を分析する方法とフローサイトメリーの原理を使った 2 つの方法がある。しかし、パターン認識技術を使用したものは、熟練した検査技師が顕微鏡を使って細胞観察を行う方法を機械化したものであり、スライドによる標本作製や一つ一つの細胞を画像処理するなど、分析に多くの手間と時間を要する。また、細胞の微細構造を人間が行うように検出することが困難であり、異常細胞の検知性能は十分とは言えず、機器も大がかりで高価な状況にある。これに対して、フローサイトメリー法は、吸光度や散乱、蛍光を用いて細胞分類する技術が発達しており、数十秒で大量の細胞を分析し検査結果が得られることから、多くの施設で利用されている。

4-1-4 白血球分類・フローサイトメリーの原理

臨床検査分野の自動血球分析装置に組み込まれている白血球分類装置は、フローサイトメリーの原理を利用したものである。自動分析装置に導入されることで高い処理能力が得られ、一次スクリーニングとして有用性が高いものである。フローサイトメリーの基本原理は、細胞浮遊液中の細胞がフローサイトメータ内に導入されると単一粒子(細胞)に整えられ、レーザー光が側方から当てられると個々の粒子(細胞)の特徴に応じてレーザー光は吸収光や散乱光となり、粒子(細胞)の性質を捉えるものである(図 4-4)。近年では、分析精度を上げるため特殊な細胞染色を行い、蛍光を発するようにしたもので登場している。

光学的処理系は、粒子(細胞)にレーザー光(633nm または 488nm)が当たると図 4-5 に示すように前方散乱光、側方散乱光や側方蛍光の光が発生し、その信号についてコンピュー

タを用いて解析する。解析結果は二次元散布図として表され、決められた領域にプロットされた粒子(細胞)の数をカウントすることで白血球分類を行うものである。例として図 4-6 に堀場製作所の自動血球分析装置 Pentra MS CPR と図 4-7 にシスメックス社の自動血球分析装置 XN-1000 からの二次元散布図を示す。Pentra MS CPR は、予め設定された境界線で区切られた部分に出現した細胞数をカウントして、リンパ球(Ly)、単球(Mo)、好中球(Neu)、好酸球(Eo)をカウントする方法に対して、XN-1000 は Mahalanobis 楕円に関連付けた領域から細胞数をカウントする仕組みとなっている。その他のメーカーについても、フローサイトメリーの原理は同じであるが、メーカーによってプロットされる細胞の位置が異なるとともに、細胞分画をカウントする計算方法も異なる。

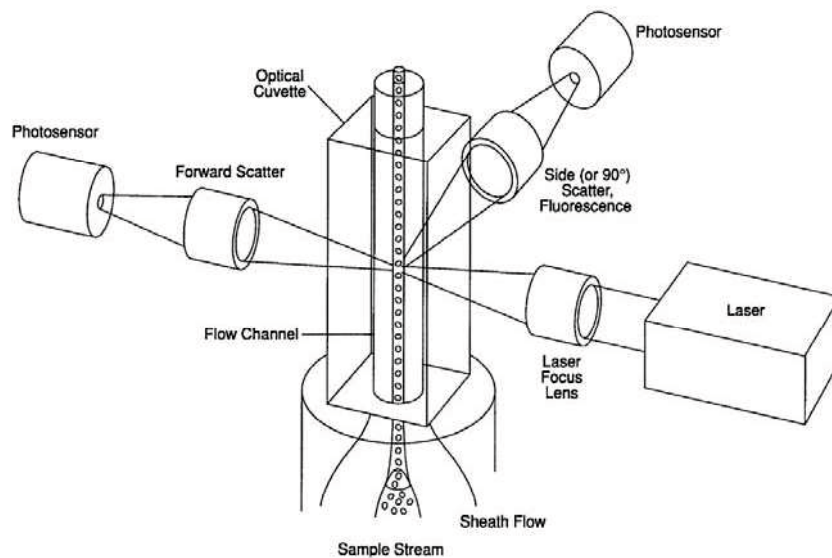


図 4-4 フローサイトメータの模式図
Differential Cell Count より

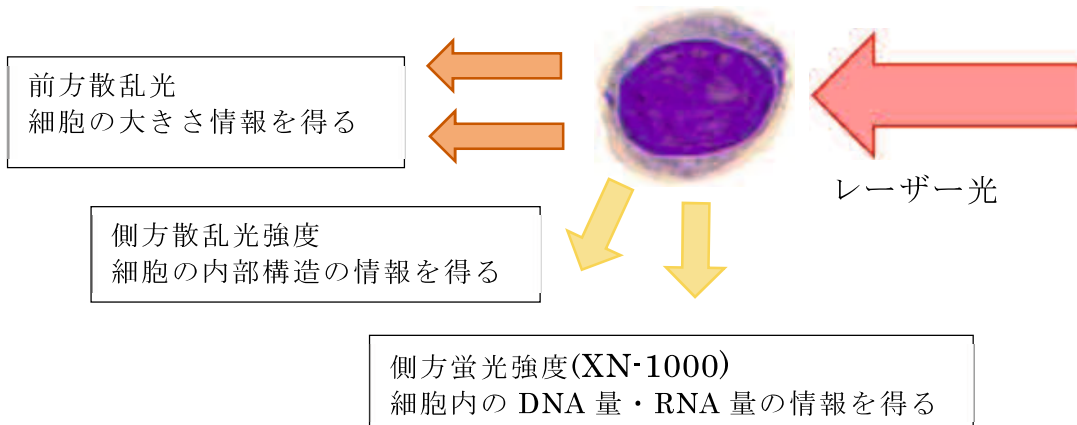


図 4-5 フローサイトメリーの光学的データ処理図

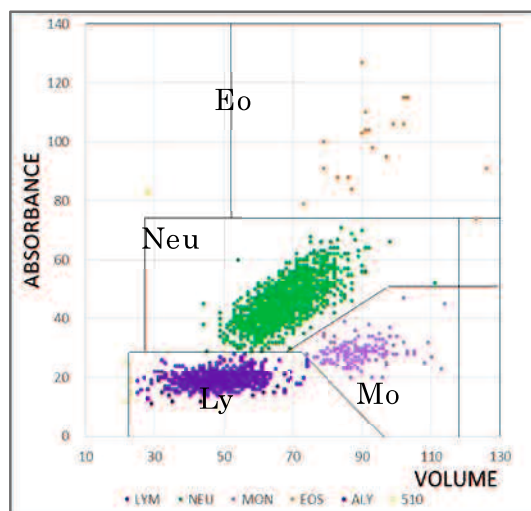


図 4-6 堀場製作所の自動血球分析装置 Pentra MS CPR からの二次元散布図
Pentra MS CPR では、直線で固定された領域に出現したデータ（細胞）数をカウントする方法をとっている。（リンパ球[Ly], 単球[Mo], 好中球[Neu], 好酸球[Eo]）

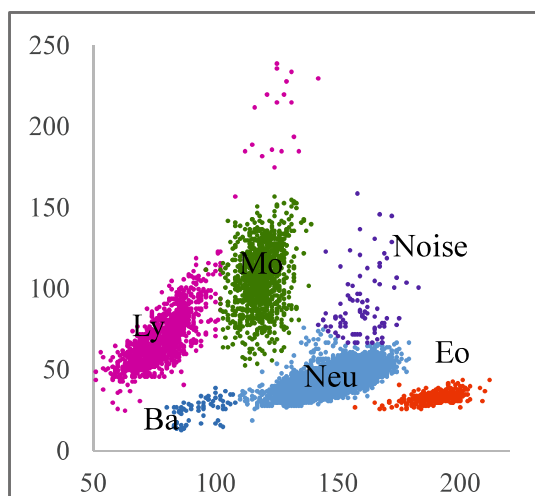


図 4-7 シスメックス社の自動血球分析装置 XN-1000 からの二次元散布図
Mahalanobis 楕円に関連付けて、その領域にあるデータ数をカウントする方法を使用している。カウント領域は、ある程度可変である。

4-1-5 白血球分類フローサイトメリー二次元散布図の特徴

堀場製作所の Pentra MS CPR とシスメックス社の XN1000 の二次元散布図が示すように、分析機器によって各細胞の出現する領域が異なっているが、好中球、単球、リンパ球は近接する位置にあり、好酸球は独立してある程度固定された位置に検出される。

これは、好酸球が好酸性の特殊顆粒を持っているために、フローサイトメリーでは他の細胞との区別がしやすいという特徴からである。その他は細胞間の特徴が近似しているために近い位置に検出され、ときに混在することもある(図 4-8)。

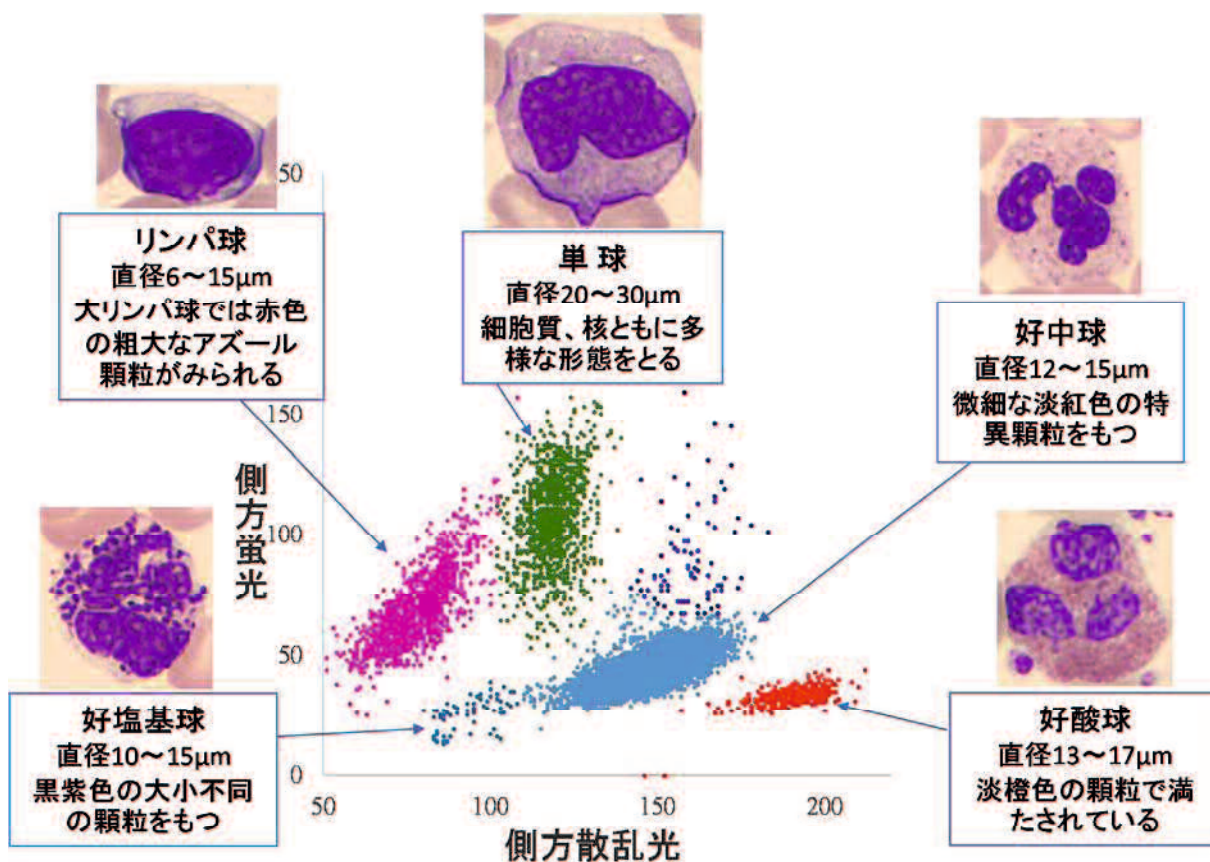


図 4-8 シスメックス社の XN-1000 の各細胞分布と細胞形態

シスメックス社の XN-1000 の二次元散布図とその位置に相当する細胞の写真を示す。好酸球は大きな特殊顆粒があるために安定的に独立したクラスターとして捉えやすいが、他の細胞は形態的に近似する部分があるために、時に非常に近接した位置に出現することがある。

4-1-6 塗抹鏡頭における白血球分画の統計的境界と手技的境界

白血球分画の基準法は目視法である。白血球分類を行うためには熟練した技術が求められる。しかし、正確なデータを出すためには、日常検査における目視する細胞数が200個程度であるため、細胞比率が少ない細胞においては統計的な限界がある。図4-9は、目視法における統計的な限界を示す。例えば単球のように細胞比率が少ない場合を考えると、細胞出現率が5%のときには200カウントの場合の二項確率から95%信頼区間は2.0~8.0%となる。400カウントの場合でも3.5~7.25%と幅が広い。これが一般的な自動血球分析装置で行うように8000カウント測定となると、信頼区間は4.5~5.5%まで狭くなる。このように、カウント数が少ない場合には統計的誤差が発生しやすい状況にある。

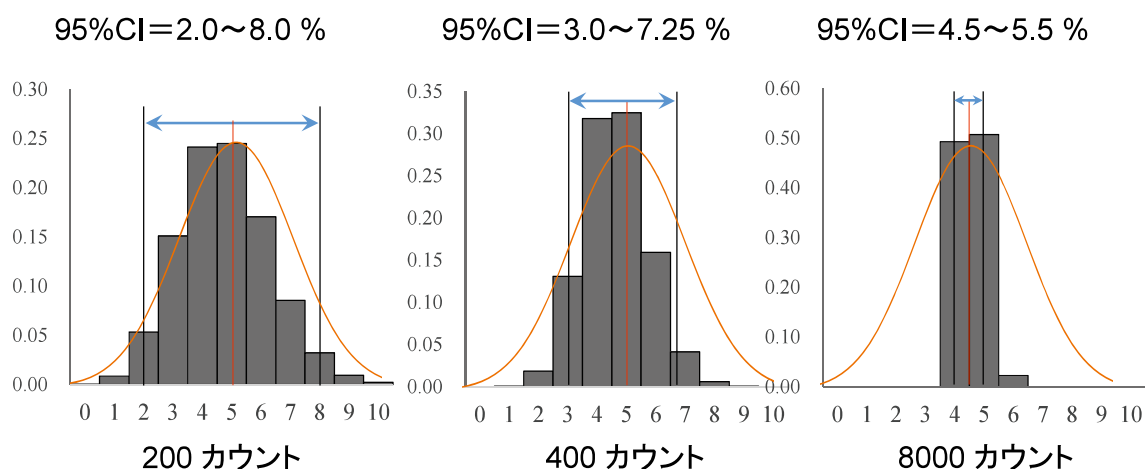


図 4-9 細胞比率 5%の 200 カウントと 400 カウントの 95%信頼区間の比較

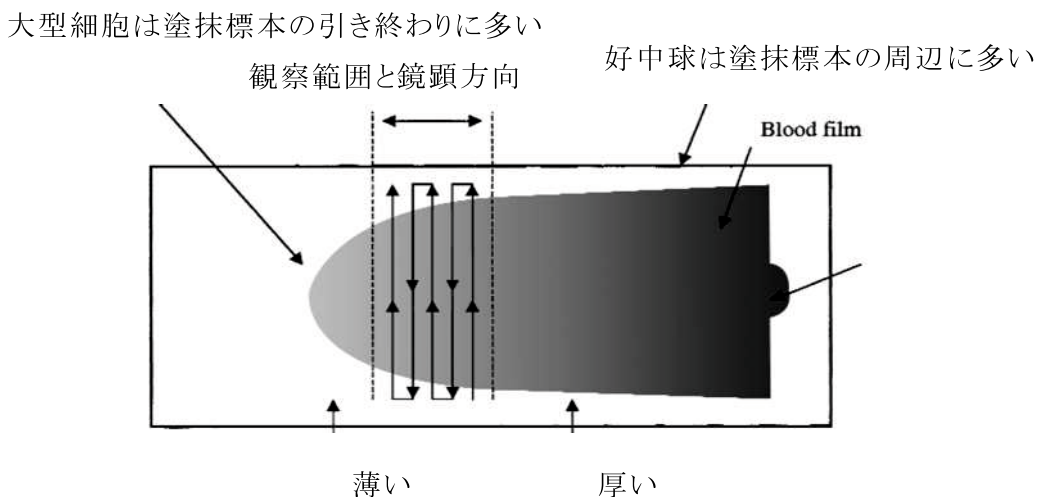


図 4-10 ウェッジ法による血液塗抹標本の状態

Houwen B The differential cell count. Lab Hematol. 2001 より

また、塗抹鏡顕では、ウェッジ法でスライドガラス上に血液を薄く塗布した血液塗抹標本を観察する(図 4-10)。しかし、血液をスライドガラス上に均一に塗布する事は困難であり、単球や好中球などの大型細胞は塗抹標本の引き終わりや周辺に多くなってしまう傾向がある³¹⁾。このような不均一な状況は、計測者の鏡顕する位置によって分画データが変わることを意味し、安定した分画データを報告することが難しい状況であることがわかる。したがって、精度の高いデータとするためには、多くのスライドを作成してより多くの細胞カウントを行う必要がある。ただしこのような作業は日常検査では困難である。

4-1-7 フローサイトメリーの限界

血球自動分析装置では、健常人検体では約 8,000 個の細胞を計測するため、その分析性能は顕微鏡的分析法(塗抹鏡顕)をはるかに凌ぐ性能を有している。しかし、細胞数の少ない単球や好酸球については 4-1-6 に示したように統計学的な誤差も多くなるため、目視法と分析機器ばかりでなく、各種の機器分析間においても相関性が悪い状況である³²⁾。Tan は目視 400 カウントと機器分析の相関をとり、好中球、好酸球やリンパ球が相関係数 0.95 以上であるのに対して、Mo の相関性は 0.75。Ba では 0.58 と低いことを報告し、このような相関性の低さは一般的に見られる現象であるとした³³⁾。

また、白血病細胞や異形リンパ球のような病的細胞では、二次元散布図に展開される細胞集団の領域が健常人とは異なるため、二次元散布図の領域を固定した設定で分類するシステムでは白血球分画性能は低下し、的確な分類を行う事が困難となる。また、目的細胞以外の死細胞や凝集物、残渣なども計測してしまうため、これらバックグラウンドノイズの影響を受けてしまい、正しい計測が行えないという問題がある。バックグラウンドノイズの影響を受けずに解析するシステムが要求される。

図 4-11 は、XN-1000 における特徴的なデータの二次元散布図とその立体図を示す。グラフと数値データから各細胞の重心位置とクラスター内分散および密度が大きく異なることが分かる。特に好中球(Neu)は細胞数が多く、かつ急峻なデータ分布であることが立体図からうかがえる。逆に単球(Mo)は、数が少なく、分布幅も広いこと事から立体図で確認することも難しいことが分かる。このように、検体ごとに大きく細胞の分布状況が異なる。

また、図 4-12 は、無作為抽出した 100 人分の各細胞のクラスター重心をプロットしたものである。グラフを見ると検体ごとに各細胞の重心位置にずれがあり、特に好中球では極端に異なる位置に細胞集団重心がある場合がある。これは、生体反応により幼弱な細胞が出現したために、このような重心位置になったと考えられる。リンパ球と単球の重心位置に範囲は、好中球ほど大きくないが上下に長く伸びた形を示しており、XN-1000 の特徴である側方蛍光の大きさに違いが発生していた³⁴⁾。現行のフローサイトメータは、細胞測定領域を固定しているものが多く、検体の個体差による細胞出現位置のずれや病態による細胞出現部位の変化に対応しきれない。個体差や病態によって細胞に出現位置の違いにも影響を受けない、細胞の動きに追従したクラスター分析法が必要と考えられる。

泳動分析

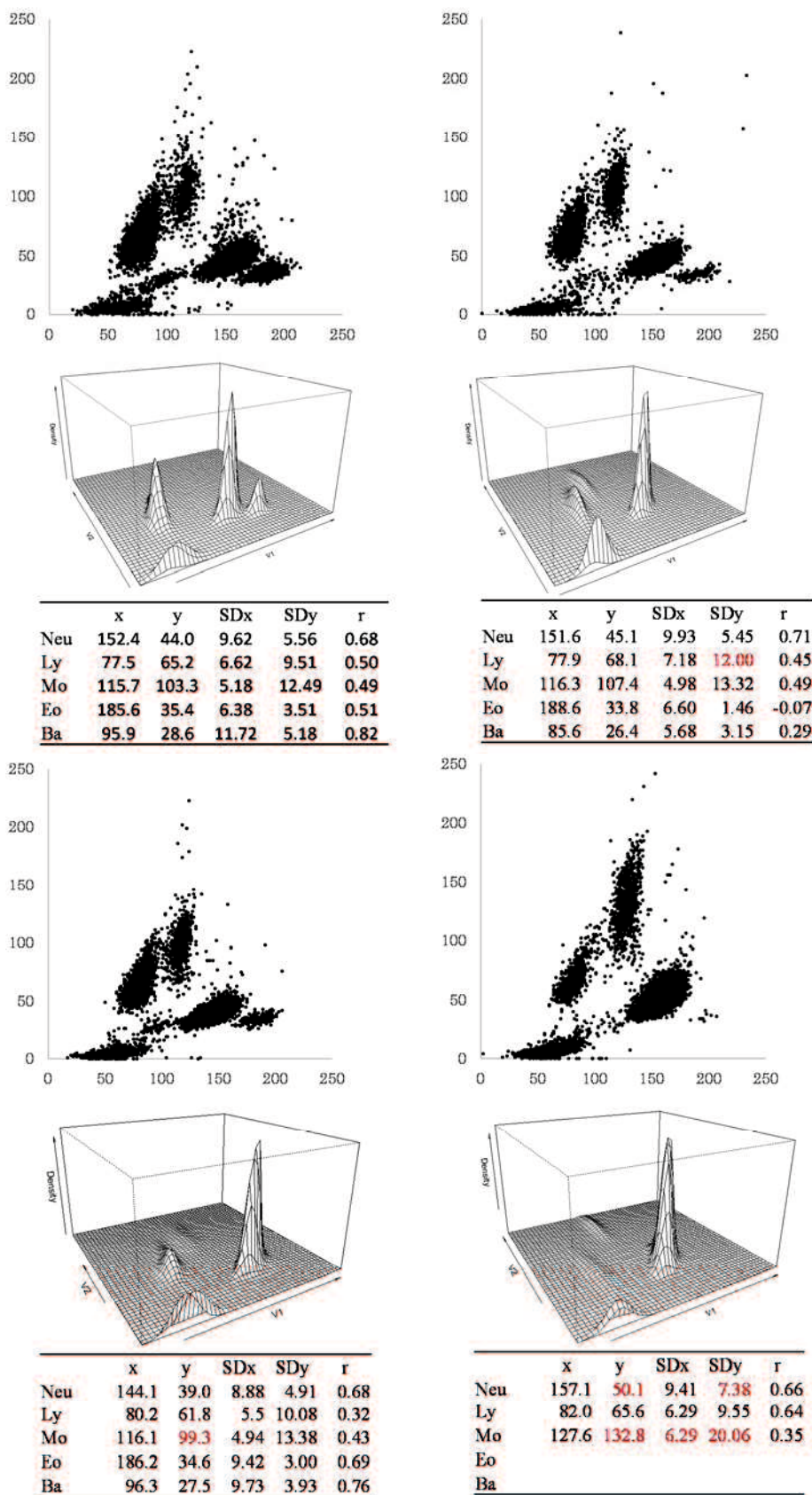


図 4-11 XN-1000 における特徴的なデータの二次元散布図

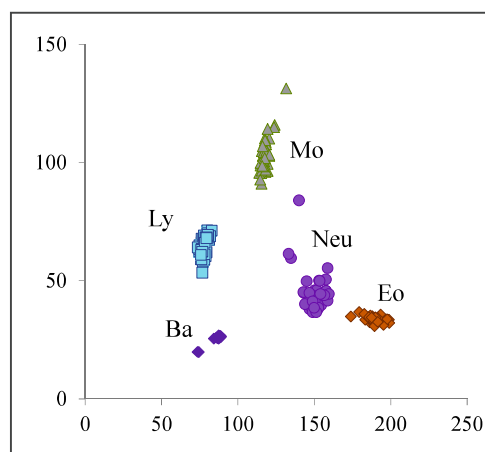


図 4-12 100 人の無作為抽出した各細胞のクラスター重心点をプロットした図
Eo と Ba はクラスター重心位置の変動は少ないが、Ly と Mo は縦に伸びた変動が認められ、Neu は極端に上部に出現する例が認められた。

4-2 実データによる検証

4-2-1 SPC/ITC 法と EM アルゴリズムの性能を評価

図 4-13 は、白血分画の代表的なデータセットとして、健常人症例 2 件と白血球減少例、白血球増多例の計 4 症例について、Pentra MS CRP によって得られた二次元散布図である。SPC/ITC 法と混合分布を適切に分画する EM アルゴリズムを比較した。SPC/ITC 法の結果は、血液検査の専門家によって得られた白血球分画データとの比較においても満足する結果が得られた。対照的に、EM アルゴリズムは適切に白血球サブタイプを識別することができず、不自然なクラスターを作成した。EM アルゴリズムでは本来 4 つのクラスターとなるところが、高密度な部分は過分画となり、ノイズの影響を受けてクラスターの存在していないところにクラスターがあるように判定がなされている。原因として①バックグラウンドノイズの存在、②高密度のクラスターに対する極小領域における局所的最適解があげられる。XN-1000 についても、別検体であるが健常人検体と白血球増多例について比較を行ったところ、SPC/ITC 法では細胞集団を上手く捉えていたが、EM アルゴリズムでは Pentra MS CRP と同様にバックグラウンドノイズや局所的最適解の影響で不適切なクラスターとなった。

現行の複数機器による検証結果から、現状のクラスター分析法では、実データには適用できないことが明らかであるが、SPC/ITC 法は臨床検査分野での活用が十分可能であると考えられた。

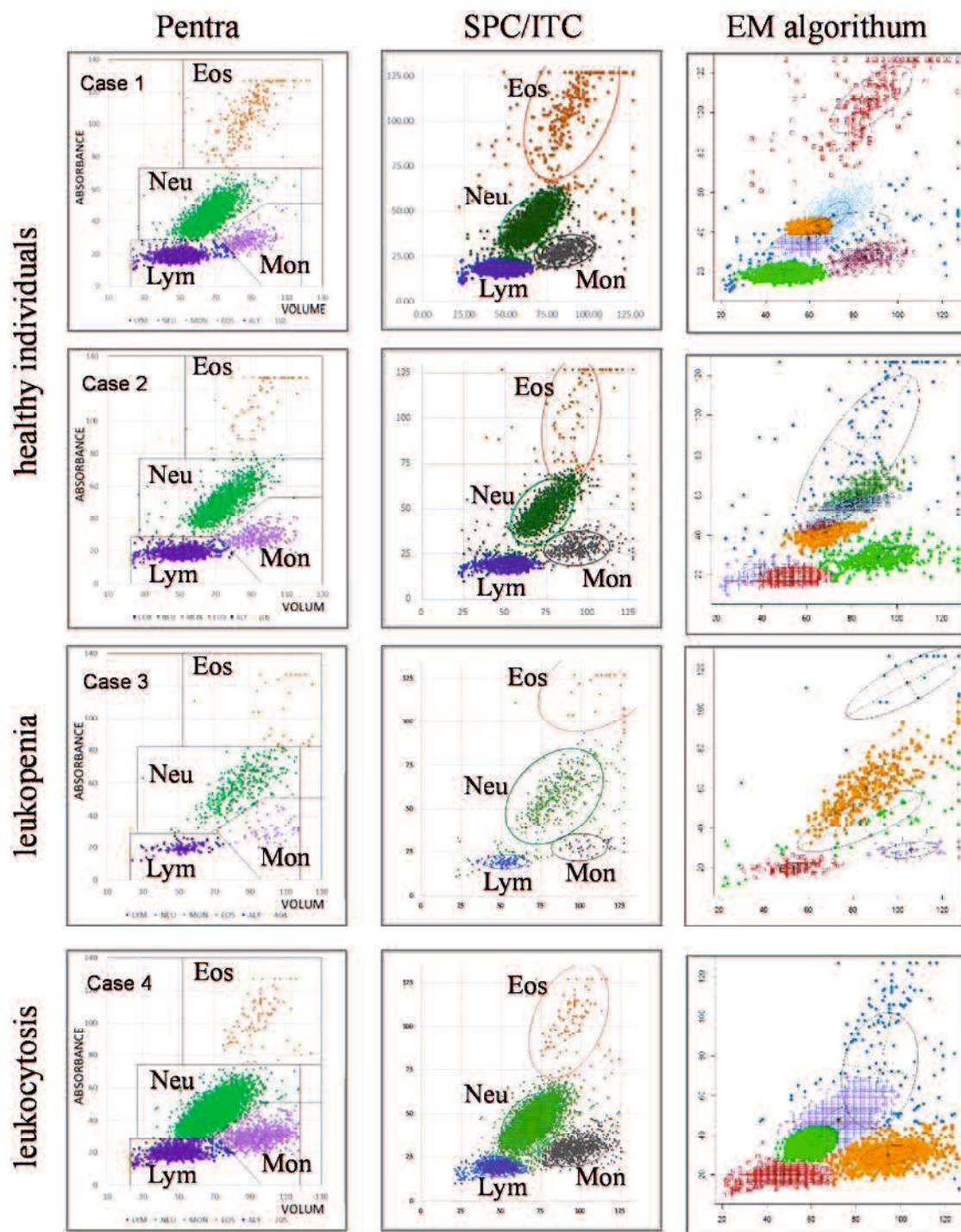


図 4-13 Pentra MS CRP における SPC/ITC 法と EM アルゴリズムの性能評価
 代表的な白血球分画の 4 つのデータセット内について SPC/ITC 法と EM アルゴリズムで比較した。EM アルゴリズム(mclust 使用)は、ノイズに対してクラスター形成が起こり、また Neu のように多数の細胞で構成されるクラスターが過分画される傾向があった。SPC/ITC 法は、すべてのケースにおいて分画に成功した。(mclust で実施した EM アルゴリズムの楕円は、50%の信頼限界であることを注意)

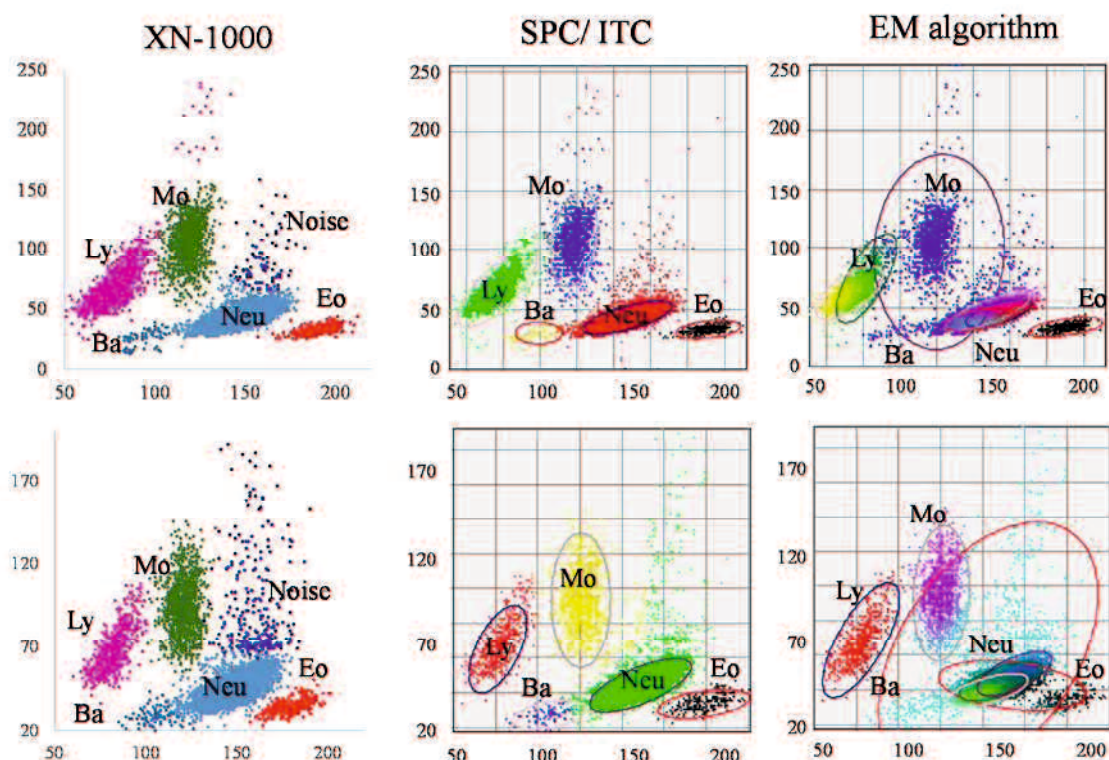


図 4-14 XN-1000 における SPC/ITC 法と EM アルゴリズムの性能評価

図 4-13 と同じく SPC/ITC 法は、適正なクラスター形成をしたが、EM アルゴリズムは細胞数が多い Neu で過分画されている。

4-2-2 健常人検体による目視法、機器分析、SPC/ITC 法の比較

京都府立医科大学病院の Pentra MS CRP によって得られた日常検体 227 例について分析を行った。検体は 1 日 10 検体のペースで抽出した 227 例の内 118 件の健常人検体を使用して SPC/ITC 法の相関分析を行った。目視法は、血液の専門家 2 名によって各 200 カウントした結果を採用した。健常検体の定義は、次の基準を満たすものとした。

$3.3 \leq \text{WBC} \leq 8.6 \times 10^9/\text{L}$, $37 \leq \text{Neu} \leq 72\%$, $20 \leq \text{Lym} \leq 50\%$, $0.6 \leq \text{Eos} \leq 8.3\%$ である。

SPC/ITC 法の基本となるデータは Pentra MS CRP アナライザから得られたものである³⁵⁾。

図 4-15 は、4 種類の白血球細胞について、目視法と SPC/ITC 法、目視法と Pentra MS CRP、SPC/ITC 法と Pentra MS CRP の相関図である。好中球における SPC/ITC 法と目視法の相関係数は $r = 0.835$ で、標準主軸回帰は $y = 0.98x + 1.57$ 、リンパ球では $r = 0.861$ 、 $y = 0.97x - 2.32$ 、好酸球では $r = 0.895$ 、 $y = 1.16x + 0.18$ で良好な結果が得られた。ただし、単球においては $r = 0.622$ 、 $y = 1.10x + 3.03$ と相関係数が他の細胞に比べて悪かった。原因は、前述の 4-1-6 に示した塗抹鏡頭における白血球分画の統計的限界と手技的限界に

よる誤差要因と単球を識別する難しさが上げられる³⁶⁾。この結果は目視法と Pentra MS CRP アナライザにおいても同様のものであった。一方、SPC/ITC 法と Pentra MS CRP アナライザの相関はどの細胞においても良好であった。これは、SPC/ITC 法の元データが Pentra MS CRP アナライザの細胞分画データを利用しているためと思われる。なお、付録として健常人および異常検体の Pentra MS CRP による二次元散布図と SPC/ITC 法による解析図を掲載した。

別の評価法として、目視法と SPC/ITC 法では、 χ^2 検定による不一致率を求めた。 χ^2 の有意水準 0.01 として、目視法と異なる比率を有する症例の割合は好中球で 4.2%(5/118)で、リンパ球は 0.8%(1/118)、好酸球 1.7%(2/118)、単球 22.9%(27/118)で、やはり単球の成績が最も悪かった。

4-2-1 異常検体による目視法と Pentra MS CRP、SPC/ITC 法の性能

227 検体中の異常検体 109 件の相関関係について図 4-16 に示す。異常検体は、血液学的悪性腫瘍、化学療法を含む検体で、健常検体以外と判断された検体とした。異常検体の計測範囲は次の通りである。 $0.8 \leq \text{WBC} \leq 31.7 \times 10^9/\text{L}$, $12.5 \leq \text{Neu} \leq 95.8\%$, $1.8 \leq \text{Lym} \leq 81.8\%$, $0.3 \leq \text{Mon} \leq 20.8\%$, $0.0 \leq \text{Eos} \leq 28.3\%$ である。

目視法と SPC/ITC 法の相関と標準主軸回帰は、Neu: $r = 0.902$, $y = 0.94x + 6.65$; Lym: $r = 0.948$, $y = 0.85x - 0.60$; Eos: $r = 0.968$, $y = 1.00x + 0.50$; and Mon: $r = 0.674$, $y = 1.39x + 1.10$ であり、目視法と Pentra MS CRP では、Neu: $r = 0.854$, $y = 0.96x - 2.88$; Lym: $r = 0.943$, $y = 0.92x + 2.48$; Eos: $r = 0.895$, $y = 0.92x + 2.48$; Mon: $r = 0.515$, $y = 2.45x - 1.82$ であった。測定範囲が健常検体よりも広がった影響もあって相関係数は総じて良くなっていた。ただし、極端に異なる症例も散見された。違いが出る症例の特徴は、好中球では低値データに集中しているが、逆に単球では高値検体で大きな違いが発生している。Buttarelli らは、白血球分画の誤差要因として塗抹の不均一、テクニック、統計的問題をあげ、最も統計的問題が大きいとしている³⁷⁾。

目視法の場合、白血球数が $2,000/\mu\text{L}$ と少ないときには、日常的細胞計測数の 200 カウントもしくは 400 カウントになかなか達しないため、本来カウントしないスライド周辺の細胞までカウントに含めてしまうことがある(スライドの周囲には大型の単球や好中球が多くなりスライド上に不均一に分布する)。このことが、目視法での誤差の原因と考えられる。細胞数が少ない状況では、目視法よりも多くの細胞を数えるフローサイトメトリーの方が正確なカウントを示すと思われる。一方、単球が高値検体でばらついている原因は、①細胞数が多い場合、反応性の血液

細胞が検出されるためにフローサイトメトリーでは分類が困難である事。②Pentra MS CRP では、好中球が多い検体の場合に、好中球の計測領域を超えて単球領域まで細胞分布が広がってしまうため、好中球を単球としてカウントしてしまう現象が起こったためと思われる。また、単球については、Pentra MS CRP と SPC/ITC 法との間においても相関性が悪かった。この原因については、4-2-2 Pentra MS CRP の細胞区画を超えた SPC/ITC 法の分析結果で解説する。

目視法と SPC/ITC 法の χ^2 検定による不一致率の評価では、好中球で 28.4%(31/109)で、リンパ球は 4.6%(5/109) ,好酸球 10.1%(11/109)、単球 30.3%(33/109)で健常検体に比較して成績が悪い。これは、計測が拡大したことや異常細胞が計測を困難にしていることが上げられ、特に目視法の限界問題の関与が考えられる。

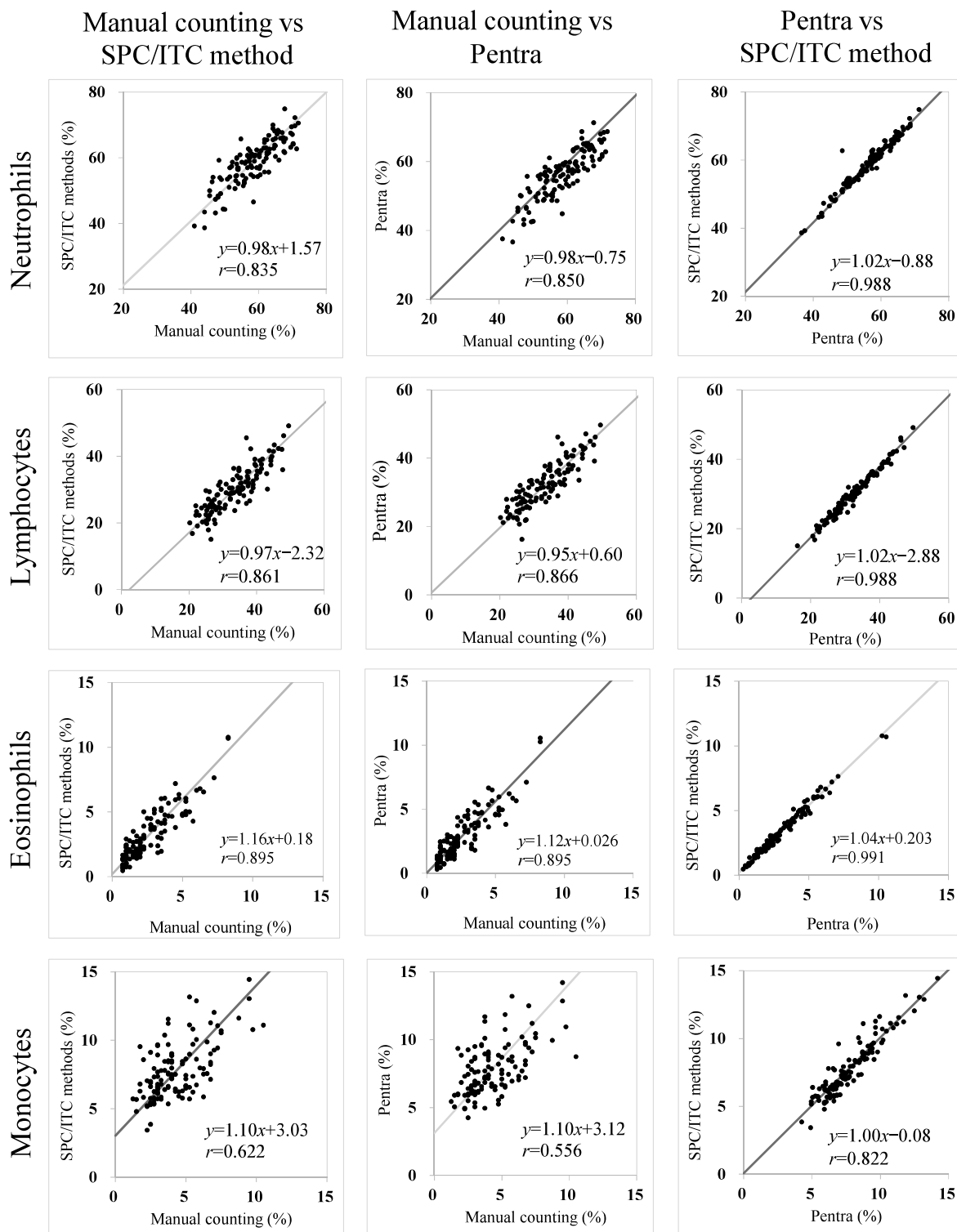


図 4-15 健常検体による目視法、Pentra MS CPR、SPC/ITC 法の比較

目視法と Pentra MS CPR は、Mo を除きおおむね良好な相関性を示した。Pentra MS CPR と SPC/ITC 法は、SPC/ITC 法が Pentra MS CPR のデータを利用して計測している関係上一致率は非常に良好であった。ただし、Mo に関しては他の細胞に比較して悪い結果であった。

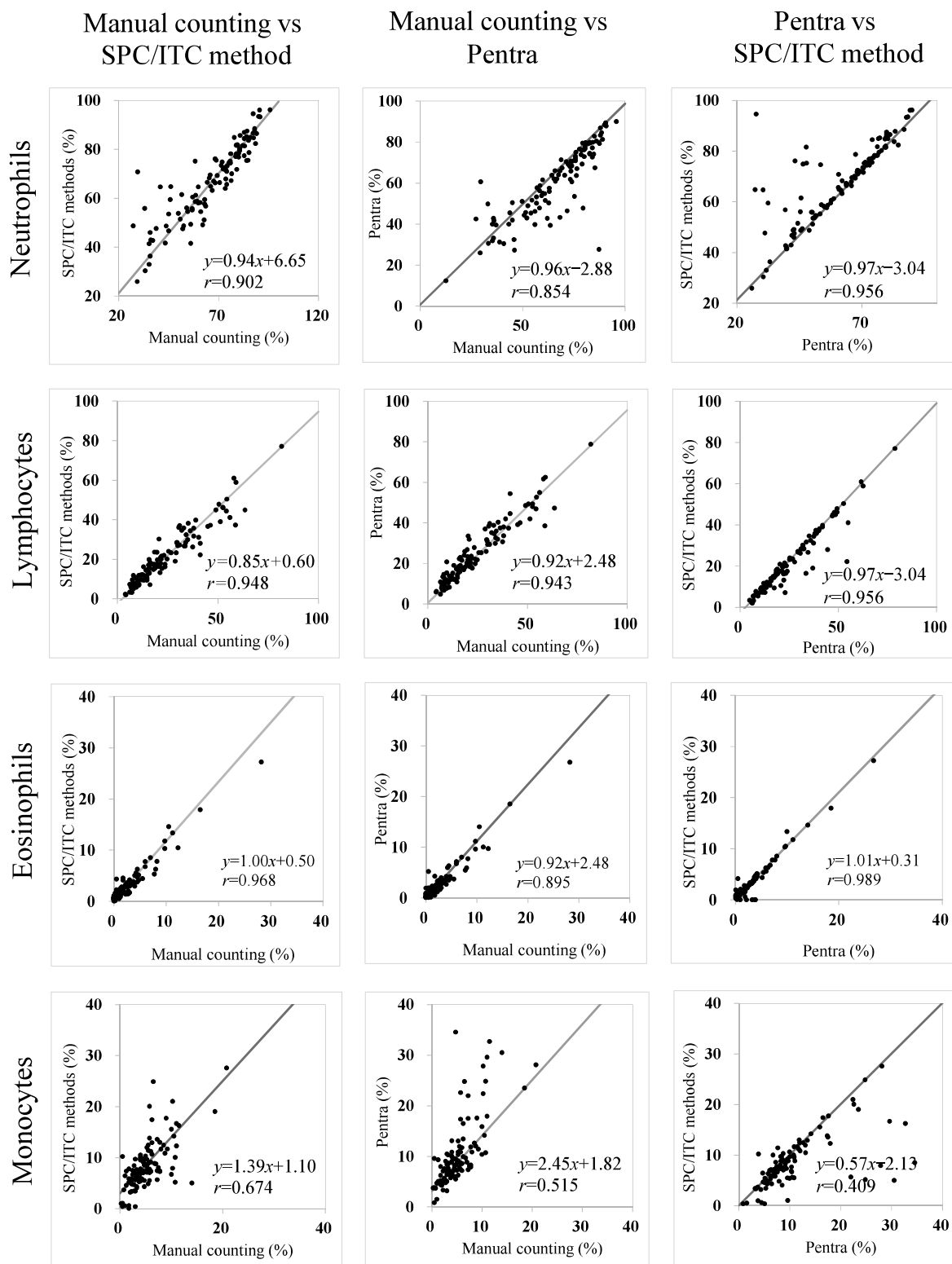


図 4-16 異常検体による目視法、機器分析、SPC/ITC法の比較

異常検体では、健常検体に比較して逸脱するデータが散見された。健常検体の Neu のときには見られなかった極端に乖離するデータが出現している。

4-2-2 Pentra MS CRP の細胞区画を超えた SPC/ITC 法の分析結果

図 4-17 は、白血球増多によって検出領域が拡大した Pentra MS CRP の細胞分画と SPC/ITC 法の分析結果を示す。パネル A から C は、好中球が単球領域まで入り込み、Pentra MS CRP で単球数を多くカウントしてしまった例を示し、パネル D は好中球が好酸球領域まで拡大して、好酸球数を多くカウントした例である。矢印(◄)は領域を超えたところを示している。好中球は炎症などの反応によって幼弱細胞が出現し、二次元散布図上で細胞出現領域が拡大しやすい傾向がある。そのため、他の細胞領域に入り込んでしまい、入り込まれた細胞は偽高値を示し、好中球は偽低値を示すことになる。この時、SPC/ITC 法では、クラスターの動きに合わせて、自在にクラスター範囲を変えるため、適切な細胞数カウントが行われていることを確認した。ただし、パネル D のように好酸球では極少数もしくは検出できない検体があり、このような場合は、SPC/ITC 法においても固定の領域のデータを参考にするなどの処置が必要となる。

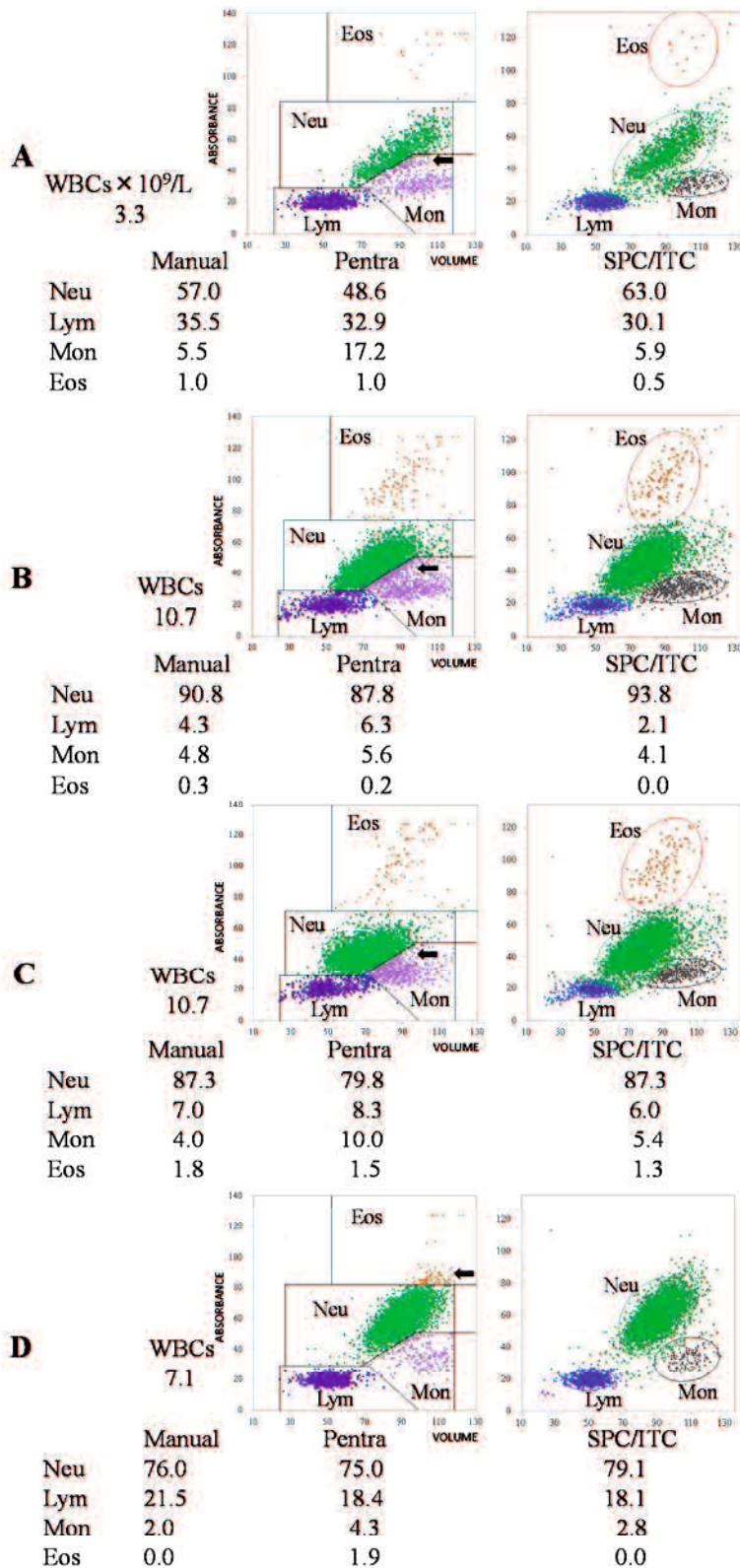


図 4-17 Pentra MS CRP の細胞区画を超えた時の SPC/ITC 法のクラスター分析結果

異常を示す代表的な 4 つのデータセットを示す。Pentra アナライザは細胞計測領域が固定であるため、境界を越えるような場合は白血球分画に問題が生じる。矢印で示す部分は、SPC/ITC 法はクラスターの変化に追従するため、矛盾なく細胞計測を実施することが可能であった。

4-3 白血球分画の実データによる各種クラスター解析のまとめ

白血球分画を機器分析で実施することは、古典的な課題を持っており、自動血液分析装置のメーカーは、様々な状況に対処するための考案が続けられている。フローサイトメトリー法での 2 次元散布図内の細胞分類は、クラスター分析を適用するための格好の場であるが、現行のクラスター分析法では、クラスター分析上の課題があり、細胞分類を実施することができていない。

2 次元散布図中の固定された境界線を使って白血球細胞分類を行う方法は、異常検体において正確でないことを説明した。新クラスター分析法は、各種の統計処理法を合成したクラスター分析法で、仮想クラスターの初期重心点の検索方法として IP 方式クラスター探索法を使用し、距離に関しては、Mahalanobis 距離を使用してデータをクラスターに属する確率として捉え、二次元反復切断補正法を使ってクラスターの中心領域にあるデータからクラスター全体を推定する。これらの方法は、オーバーラップや極端値、ノイズ、極端なデータバランスの違いに強い方法である。

実データでは、各クラスターが正規分布を呈しておらず、バックグラウンドノイズ成分も多く含まれ、かつクラスター間のデータ数も大きく異なるものが出現する。複雑な混合分布を的確に処理する EM アルゴリズムであってもノイズの影響を大きく受けてしまい、ノイズを一つのクラスターと見なして処理が行われている。このような実データにおいても SPC/ITC 法によるクラスター分析は適切に動作しており、ノイズの影響を受けずに十分な性能を発揮して、細胞を的確に捉えた。

また実データにおいては、血液細胞数が大きく異なり、好中球は健常人で 40~70%を占めるが、生体内で炎症反応が起こるとその数は 10,000/ μ L 以上にも及ぶ。このような検体で EM アルゴリズムを適応すると極小領域で局所最適解を検出してしまうために過剰なクラスター形成が起こる現象が認められた。このような検体に対しても SPC/ITC 法は適切なクラスター形成ができ、その有効性を確認した。本法は、実データに対して頑健性の高い方法といえる。

一方、標準法である目視法と機器による分析法については、目視法の統計的問題(200~400 カウントのため、少ない細胞種の場合に誤差が大きい)や大型細胞特有のスライドガラス上の分布の偏りのため、単球のような検出率が少なく、大型細胞における正確性には問題がある。また、単球は熟練技術者同士の目視法であっても一致率が低い細胞であるため、相関性が悪いことが指摘されている。機器分析は細胞のカウント数は目視法の 10 倍以上を計測していることから、統計的な誤差は小さく、細胞の分布に関しても均一な状態で計測が行えるため、原理的に単球や好酸球などの少ない血液細胞での計測精度は高いと思われる。

しかし、固定領域によって区別する自動血液分析装置の結果には注意が必要な場合がある。近接する集団の中には、異なる集団の中に分類される細胞がある。特に単球のカウントには、未成熟な単球が異型リンパ球としてカウントされること、好中球の細胞が増多したことによる細胞区画を乗り越えて単球領域に入ってくる現象などの様々要因があり、異常検体での計測には十分な注意が必要である。このような検体に対して **SPC/ITC** 法は、一つの細胞がどの集団に属すかを確率論的に分類する方法であるため、細胞が近接した状態や検出位置が大きく変化した場合でも適切な分類が可能である。自動血液分析装置のフローサイトメトリーを用いた白血球分類には、より適切なクラスター分析が実施可能な本法の採用が望まれる。

また、実装した **SPC/ITC** 法から得られた各クラスター平均、標準偏差、相関係数などのパラメータは、白血球分類ばかりではなく、得られる重心、標準偏差、相関係数などのパラメータは、異常検体の検出や病態解析への応用が可能である。

第 5 章

総 括

- 5-1 本研究の成果
- 5-2 本研究の限界
- 5-3 本研究の今後の展開

第 5 章. 総 括

5-1 本研究の成果

第 1 章では、クラスター分析の意味と必要性について述べた。また、臨床検査の分野での利用価値が高い方法であるが、現在のクラスター分析手法では解決できない問題があるため全く使用されていない現状である事を述べた。

第 2 章では、現行の各種のクラスター分析法について解説を行った。臨床検査の分野で有用性があると考えられる非階層型クラスター分析の欠点を取り上げ、①計算のはじめの段階で、クラスター数が決定されていなければならない、②ランダムに決定される初期値が最終結果に大きく影響する、③極小領域での局所最適解に陥りやすい、④オーバーラップするクラスターを適切に分けることが出来ない、⑤極端値やノイズの影響を大きく受けってしまうことなどの問題を示した。現在でもいろいろな場面で使用されるクラスター分析は、前述の問題点を回避するために新しい方法が考案されているが、未だ初期値に関する問題や極端値、ノイズデータに関する問題は完全に解決されていない。このような状況に対して、具体例を示して解説した。

第 3 章では、新しいクラスター分析法を考案した事を示した。現行のクラスター分析法の問題点の解決策として 3 つのアルゴリズム (IP 方式クラスター探索法、SPC 法、ITC 法) を合成した方法を解説した。IP 方式クラスター探索法は、第 2 章で示した①②の初期値問題の解決として提案した方法で、画像処理技術を使って適切な初期値を決定する方法である。また、SPC 法と ITC 法は、個々のデータに対して Mahalanobis 距離に基づく各クラスターへの帰属確率とクラスターのデータ数による重み付け、さらには二次元反復切断という特殊データ処理によって、現行のクラスター分析法の問題点③④⑤を解決したことを現行のクラスター分析法との比較を交えて示した。特にノイズ成分の多いデータに対して頑健性の高い方法であると思われる。検証のために臨床分野で良く使用される白血球分画についてシミュレーションモデルを作成し、具体的に有用性の確認を行った。

第 4 章では、新しいクラスター分析法の臨床分野への利用例について示した。具体的には、フローサイトメトリー法で行われる白血球分類を例に挙げ、分析装置ごとに特徴のある測定技術が用いられ、出力される結果も異なり、また個人や病態によって二次元散布図の形が異なる事を示した。機器による分画で困難な事は、細胞形態が近似する単球や好中球数をカウントすることや固定された二次元散布図での分類法を上げ、目視法においてはカウントする数が少ないことによる統計的な問題やスライドグラス上に不均一に分布する細胞をカウントする問題があることを示した。このような状況の中で適切な二次元散布図が得られれば、SPC/ITC 法を駆

使して適切なクラスター分析が実施できる事を実証した。このことは、臨床検査の自動化技術に大きく貢献するものと考えられた。

5-2 本研究の限界

IP 方式クラスター探索法の限界としては、重心を検索するために散布図を適当に分割するが、分割間隔によって相対密度も変化するために、最終結果に影響してしまう事が上げられる(図 5-1)。ただし、初期に実クラスターよりも多くのクラスターが検出されても、SPC/ITC 法のマージ機能によって吸収されて、最終結果はほぼ同じものとなる。解決策としては、臨床データであれば過去の事例などの積み重ねてきたデータから、適切な分割値が推定可能であることや白血球分類のようにある程度クラスターが発生する位置が特定されている場合については、適当な初期パラメータを決めて実施することも可能と思われる。

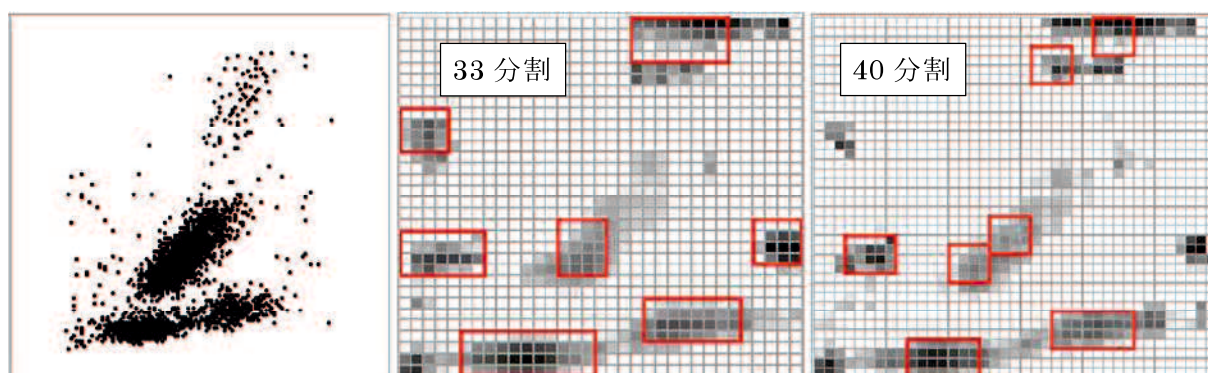


図 5-1 分割数に違いによる RD 行列の違いと初期クラスター検出結果の違い

左のパネルは、Pentra MS CRP にて測定した健常人データである。このデータについて分画数を変えた場合の結果を中央と右のパネルに示す。最初に設定する分割数によって RD 行列に違いが発生するために初期クラスターが異なる。ただし、SPC/ITC 法はマージ機能があるため、極端な分割でない限り最終結果はほぼ同じものとなる。

SPC/ITC 法の限界としては、他のクラスター分析法も同様であるが、データ数が非常に少ない場合およびバラツキが大きいデータに関しては、クラスター化が難しく、計算が困難になってしまうことである。また、クラスターが極端に近接している場合も同様である。ただし、本法は他のクラスター分析法に比較して、より近接した場合であっても適切なクラスター分析が可能であり、白血球分画の実データで示したように高度の異常検体でなければ適切なクラスター分析が可能である。白血球自体が各分画の特徴をなくしてしまっているような病的状態においては、機器による分画の限界を超えていると考えられることから、このような場合には目視法によって確認するべきと考えられる。

5-3 本研究の今後の展開

5-3-1 SPC/ITC 法から得られたパラメータの利用

SPC/ITC 法を実行するとクラスター分析結果ばかりではなく、各クラスターの重心点(平均)クラスター内分散、相関係数 r 、回帰係数 $y=a+bx$ が算出される。これを白血球分画に利用した場合、感染症などの炎症疾患時には白血球数の変動のみならず、その分画の位置や分散にも変化をきたすことから、炎症の有無やより細かく炎症の状態の把握と治療効果の判定に利用することが考えられる。また、特定の血液疾患ではその疾患に特異なパラメータを示すことが推定されることから、今後検討すべき内容と思われる。

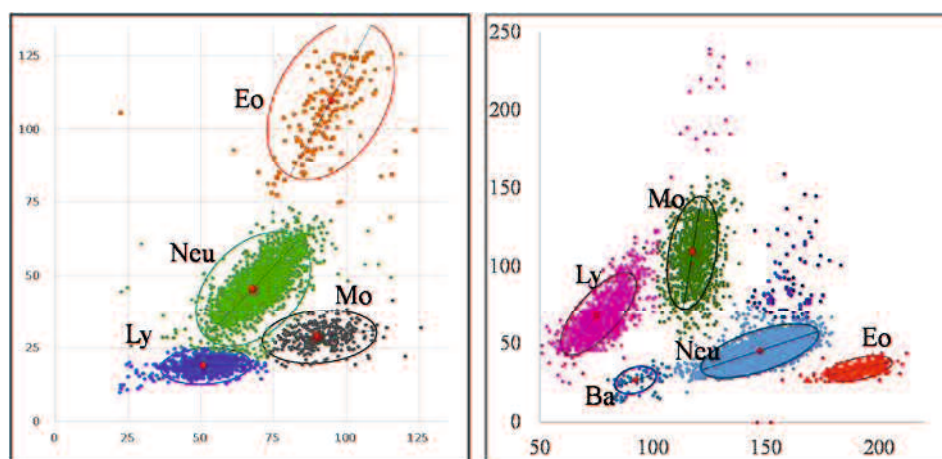


図 5-2 SPC/ITC 法から得られたパラメータ(標準偏差、相関係数、回帰式)の利用
赤い点は、各クラスターの重心点を表し、楕円はクラスターの 95%信頼区間を示す、そして、楕円の長軸に一致した直線は、標準主軸回帰である。

5-3-2 SPC/ITC 法の白血球分画以外の臨床分野への応用

SPC/ITC 法の考え方は、白血球分画ばかりではなく、1次元の混合分布となっているタンパク電気泳動などの分画にも応用可能である、また、リンパ球などの単一細胞をより詳細に分類するフローサイトメトリー法では、分画を用手法で行っているが、この分野においても SPC/ITC 法を利用した自動化に貢献できると考えられる。

5-3-3 疾患分類への応用

SPC/ITC 法の自己分配という考え方は、血球の自動分類に限定されるものではない。例えば、自己免疫疾患などの互いにオーバーラップした特異疾患群の分類にも適用可能である。自己分配方式の威力は、各症例が分離された個々の自己免疫疾患の要素をどの程度併せ持っているかを、“自己分配係数”から推定でき点である。すなわち、どこかに帰属するだけでなく、複数の病態にまたがった存在であることを推定することができる。このことは、疾病分類の新しい手法としての可能性もあり、さらなる研究が求められる。

5-3-4 統計分野への利用

SPC/ITC 法はクラスター分析そのものであるため、臨床面のみならず広く統計学分野に利用可能である。

倫理

本研究の白血球分類データは、千葉県救急医療センターの倫理委員会での承認(H27-5)および京都府立医科大学での承認(RBMR-C-1144-2)を得て実施した。

謝辞

本研究は、著者が山口大学大学院医学系研究科後期博士課程在学中に、同大学の市原清志教授の指導のもとに行いました。市原清志教授に深謝いたします。また、血液分析装置(PENTRA MS CRP)の2次元データセットの提供にご協力頂いた堀場製作所の齊藤憲祐氏および研究所の方々に感謝いたします。さらに、有益な議論をして頂いた市原研究室の皆様にも感謝いたします。

1. 杉山将. イラストで学ぶ 機械学習 最小二乗法による識別モデル学習を中心に. 講談社. 2013: 2-4.
2. Forgy E. Cluster analysis of multivariate data: efficiency vs. interpretability of classification. *Biometrics* 1965; 21:768-9.
3. MacQueen J, Some Methods for Classification and Analysis of Multivariate Observations, in: Proc. of the 5th Berkeley Symposium on Mathematical Statistics and Probability 1967; 281-97.
4. Memarsadeghi N, Mount DM, Netanyahu NS, Moigne J. A fast implementation of the ISODATA clustering algorithm. *Int J Comput Geometry Appl* 2007; 17: 71-103.
5. Arthur D. "K-means++: The advantages of careful seeding", Proc. of the eighteenth annual ACM-SIAM symposium on Discrete algorithm 2007; 1027-35.
6. 小野田崇. K-means 法の様々な初期値設定によるクラスター分析結果の実験的比較. The 25th Annual Conference of the Japanese Society for Artificial Intelligence. 2011.
7. 豊田秀樹. 変数間の関係性を考慮してクラスタ数を決定する K-means 法の改良. 32 心理学研究 2011; 82:32-40
8. Maesschalck RD, Rimboud DJ, Massart DL. The Mahalanobis distance. *Chemom Intell Lab Syst.* 2000; 50:1-18.
9. Girolami M. Mercer kernel based clustering in feature space. *IEEE Trans. on Neural Networks* 2002; 13:780-84.
10. Bezdek JC. *Pattern Recognition with Fuzzy Objective Function Algorithms.* Plenum press. 1981; 65-70
11. Liu HC, Jeng BC, Yih JM, Yu YK. Fuzzy C-Means Algorithm Based on Standard Mahalanobis Distances. *Inte Symp Inf Processing* 2009. 422-27.
12. Mohamed Jafar, Sivakuma R. Data Clustering Based on Hybrid of Fuzzy and Swarm Intelligence Algorithm Using Euclidean and Non-Euclidean Distance Metrics. A Comparative Study. *Journal of Theoretical and Applied Information Technology* 2014; 69: 599-610.
13. 八木隆文, 市橋秀友, 本多克宏. 繰り返し重み付最小2乗法と Mahalanobis 距離による FCM クラスター分析. 21st Fuzzy System Symposium 2005; 29-34.
14. Dempster AP, Laird NM, Rubin D. Maximum likelihood from incomplete data via the EM algorithm. *J Royal Stat Soc B* 1977; 39:1-38.

15. Wu CFJ. On the convergence properties of the EM algorithm. *Annals of Statistics* 1983; 11:95–103.
16. Naim I, Gildea D. Convergence of the EM Algorithm for Gaussian Mixtures with Unbalanced Mixing Coefficients. *The 29th International Conference on Machine Learning*, UK. 2012.
17. Pelleg D, Moore A. X-means: Extending K-means with Efficient Estimation of the Number of Clusters. *Proc. Of the 17th International Conference on Machine Learning* 2000: 727–34.
18. Hourdakis N, Argyriou M, Petrakis EG. Hierarchical clustering in medical document collections: the BIC-means method. *Journal of Digital Information Management* 2010; 8:71–7.
19. Fraley C, Raftery AE, Murphy TB, Scrucca L. *mclust* Version 4 for R: Normal Mixture Modeling for Model-Based Clustering, Classification, and Density Estimation. 2012. 3 Aug 2015
<<http://www.stat.washington.edu/research/reports/2012/tr597.pdf>>
20. Notsu A, Komori O, Eguchi S. Spontaneous Clustering via Minimum Gamma-divergence. *Neural Computation* 2014; 26:421–48.
21. Katsavounidis I, Kuo C-CJ, Zhang Z, A New Initialization Technique for Generalized Lloyd Iteration. *IEEE Signal Processing Letters* 1994; 10:144–46.
22. Celebi ME, Kingravi HK. Vela: A Comparative Study of Efficient Initialization Methods for the K-means Clustering Algorithm. *Expert Systems with Applications* 2013; 40: 200–10.
23. 宮崎 文吾, 和泉 潔, 鳥海 不二夫, 高橋 諒. 混合ガウスモデルを用いた市場注文状況の変化の検出, JPX ワーキングペーパー. 2013; 3:1–32.
24. Archambeau C, Lee JA, Verleysen M. On convergence problems of the EM algorithm for finite Gaussian mixtures. *ESANN'2003 proceedings - European Symposium on Artificial Neural Networks Bruges (Belgium)*, ISBN 2-930307-03-X, pp. 99–106.
25. 白井敏明: 正常値計算法:反復切断補正法における切断係数の選択. *臨床病理* 1981; 29:319–22.
26. 市原清志: 二次元反復切断補正法の考案と外部精度管理調査への応用. *臨床病理* 2001; 49:136.
27. Ichihara K, Kawai T. An iterative method for improved estimation of the mean of peer-group distributions in proficiency testing. *Clin Chem Lab Med* 2005; 43:412–21.

28. 五利江重昭 . MS-Excel を用いた混合正規分布のパラメータ推定 . SUIZANZOSHOKU. 2002; 50:243-49.
29. Redner RA and Walker HF. Mixture densities, maximum likelihood and the EM algorithm. *Siam Review* 1984; 26:195-239.
30. 金井正光編. 臨床検査法提要 改訂第 33 版. 金原出版. 2005; 406-11.
31. Houwen B. The differential cell count. *Lab Hematol.* 2001; 7:89-100.
32. Meintker L, Ringwald J, Rauh M, Krause SW. Comparison of automated differential blood cell counts from Abbott Sapphire, Siemens Advia 120, Beckman Coulter DxH 800, and Sysmex XE-2100 in normal and pathologic samples. *Am J Clin Pathol* 2013; 139:641-50.
33. Tan BT, Nava AJ, George TI. Evaluation of the Beckman Coulter UniCel DxH 800, Beckman Coulter LH 780, and Abbott Diagnostics Cell-Dyn Sapphire hematology analyzers on adult specimens in a tertiary care hospital. *Am J Clin Pathol* 2011; 135:939-51.
34. Kawauchi S, Takagi Y, Kono M, Wada A and Morikawa T. Comparison of the Leukocyte differentiation Scattergrams Between the XN-Series and the XE-Series of Hematology Analyzers. *Sysmex Journal International.* 2014; 24:1-8.
35. Inaba T, Nomura N, Ishizuka K, Yoshioka K, Takahashi, M, Yuasa S, et al. Basic evaluation of Pentra MS CRP, a new automated hematology analyzer for rapid 5-part WBC differential and CRP using a small volume of whole blood. *Int J Lab Hematol* 2015; 37:208-16.
36. Goossens W, Hovr LV, Verwilghen RL. Monocyte counting: Discrepancies in results obtained with different automated instruments. *J Clin Pathol* 1991; 44:224-27
37. Buttarello M, Plebani M, Automated Blood Cell Counts. *Hematopathology.* *Am J Clin Pathol.* 2008; 130:104-16.

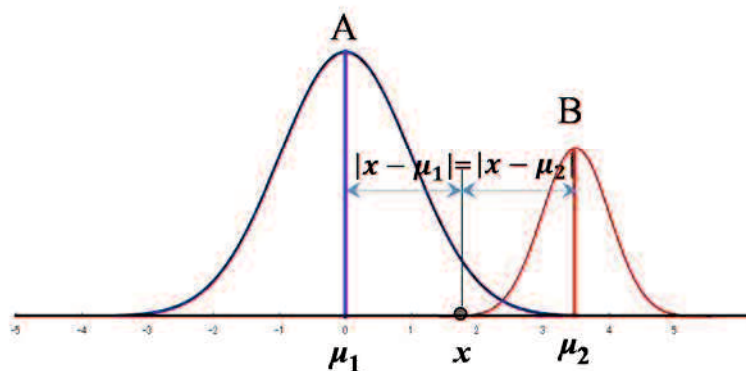
付録

● Mahalanobis 距離について

正規分布のもとで最も自然に導かれるのが Mahalanobis 距離である。集団のデータは、ある大きさとバラツキをもつものであることから、単純な距離で近さを判断することは誤りである。下図に示すような 1 次元図の場合の μ_1 と μ_2 の間に位置 x の距離について所属の面から考えてみると、 μ_1 と x のユークリッド距離(1)(単純な長さを求める式)では、 μ_1 からの x と μ_2 からの x は等距離にあり、どちらの集団に属するかは五分五分である。一方、Mahalanobis 距離(2)は分散の要素を計算式に入れて、(1)の式の二乗を分散で割った値を採るので、どちらの集団に属する確率を見ると、Mahalanobis 距離では、分散の大きい集団 A に属するか確率が大きく計算される事が分かる。Mahalanobis 距離として見た場合は、 x と集団 A の距離が短いこととなる。

$$1 \text{ 次元} \quad \text{ユークリッド距離} \quad |x - \mu_1| \quad (1)$$

$$\text{Mahalanobis 距離} \quad D_1^2 = \frac{(x - \mu_1)^2}{\sigma_1^2} \quad (2)$$



2 次元の場合のユークリッド距離は、(3)式で表され、Mahalanobis 距離は(4)式で表される。下図は、2 変間に相関性がある 2 セット A, B を散布図モデルとして表し、楕円は Mahalanobis 距離の 95%信頼区間を示す。観察された黒い点(x_0, y_0)は、青の点と赤の点のユークリッド距離で等距離にあるデータであるが、等確率楕円の形状から、青の集団に含まれるのが自然と考えられる。2 次元の Mahalanobis 距離は、データの分散 s^2 と相関係数 r を含んだ式であり、 D を二乗した D^2 は自由度 2 の χ^2 分布に従い、各点に対して所属する確率として捉えることが出来る。し

たがって、相関関係のある 2 変量に関する処理を行う場合には、Mahalanobis 距離を用いることが妥当と言える。

2 次元

Euclidean 距離

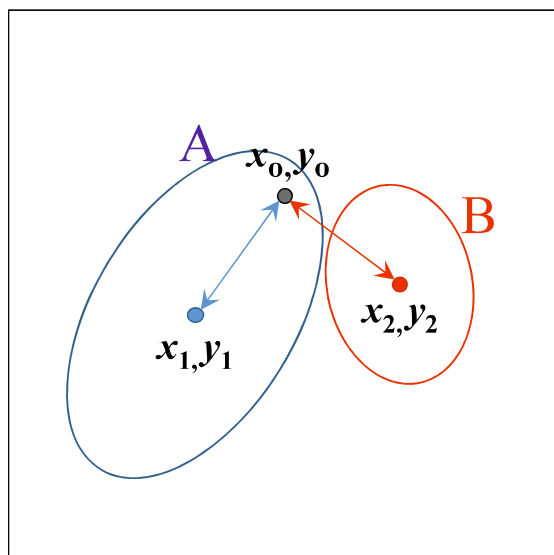
$$d = \sqrt{(x_2 - x_1)^2 + (y_2 - y_1)^2} \quad (3)$$

Mahalanobis 距離

$$D = \sqrt{\left(\frac{x_0 - \bar{x}_1}{\sigma_x}\right)^2 + \left[\left\{\left(\frac{y_0 - \bar{y}_1}{\sigma_y}\right) - \rho_{12} \left(\frac{x_0 - \bar{x}_1}{\sigma_x}\right)\right\} \frac{1}{\sqrt{1 - \rho^2}}\right]^2} \quad (4)$$

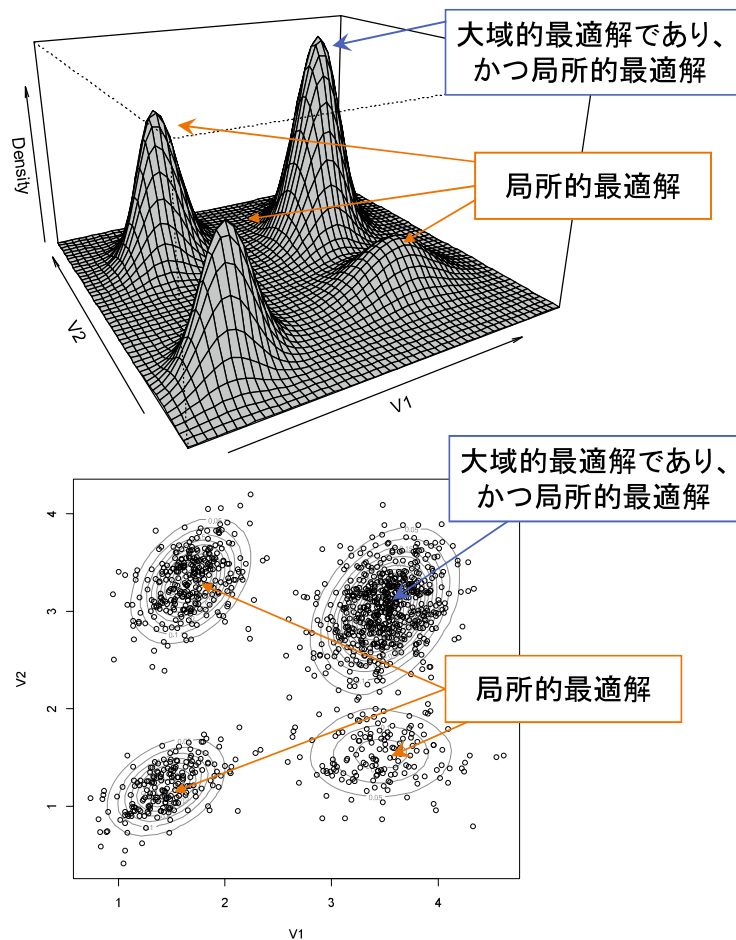
$$D = \sqrt{\frac{u_1^2 + u_2^2 - 2\rho u_1 u_2}{1 - \rho^2}} \quad (5)$$

$$u_1 = \frac{x_0 - x_1}{\sigma_x}, \quad u_2 = \frac{y_0 - y_1}{\sigma_y}$$



● 大域的最適解と局所的最適解

図は、2次元データにおける大域的最適解と局所的最適解示したものである。局所的最適解とは、限られた領域における最大値(最小値)を示す点を意味し、対して大域的最適解は全領域における最大値(最小値)を意味する。EM アルゴリズムにおいて局所最適解に陥りやすいとは、極小領域における最大値(最小値)を最適解としてしまうことである。EM アルゴリズムでは、データ数が多い場合に特に発生しやすい現象である。



● **ベイズアン情報量規準 (BIC: Bayesian Information Criterion)**

BIC は絶対値を持って評価するものではなく、比較評価をするために使用するものであり、一般的にはいろいろなモデル中から最良と考えられるモデルを選択する場合に使用する。

従来の方法では、適合度だけを評価するものであったが、パラメータ数 (モデルの説明変数) もモデルの良さを判定する基準にするものである。これによって、パラメータ数が多ければ多いほど適合度が上がってしまい、より複雑なモデルが選択されてしまうというのを防ぐ効果がある。

適合度については、どのような分布を事前分布として想定しておくかが重要なポイントとなる。しかし、最適な事前分布を決定する方法は現在の所ない。

$$BIC = -2\{\text{対数尤度}\} + \log(\{\text{サンプルサイズ}\})\{\text{モデルパラメータ数}\}$$

1. AIC (Akaike's information criterion)

$$AIC := \ln\left(\frac{S_e^2}{n}\right) + \frac{2K}{n}$$

2. BIC (Schwartz's Bayesian information criterion)

$$BIC := \ln\left(\frac{S_e^2}{n}\right) + \frac{K \ln n}{n}$$

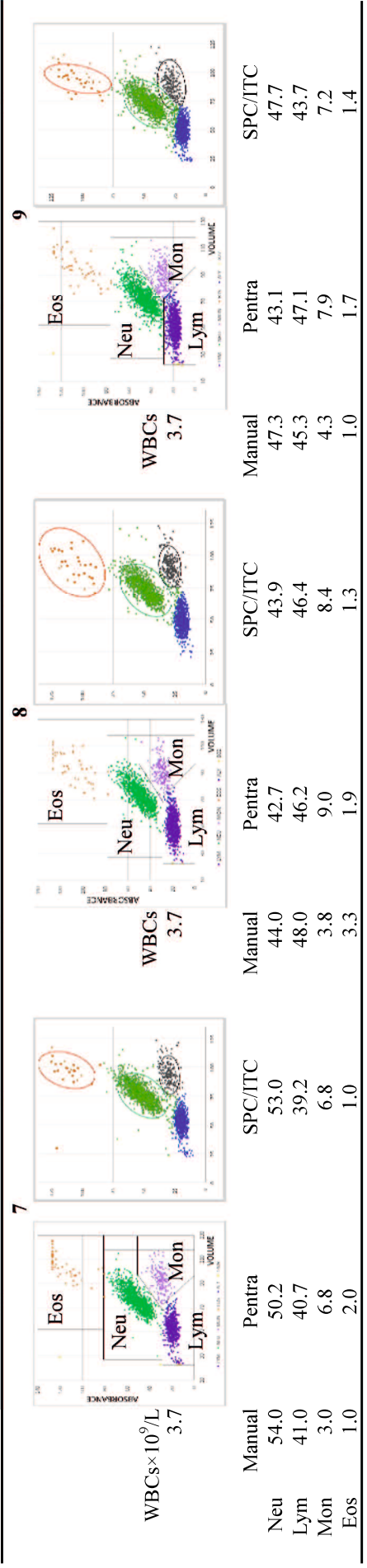
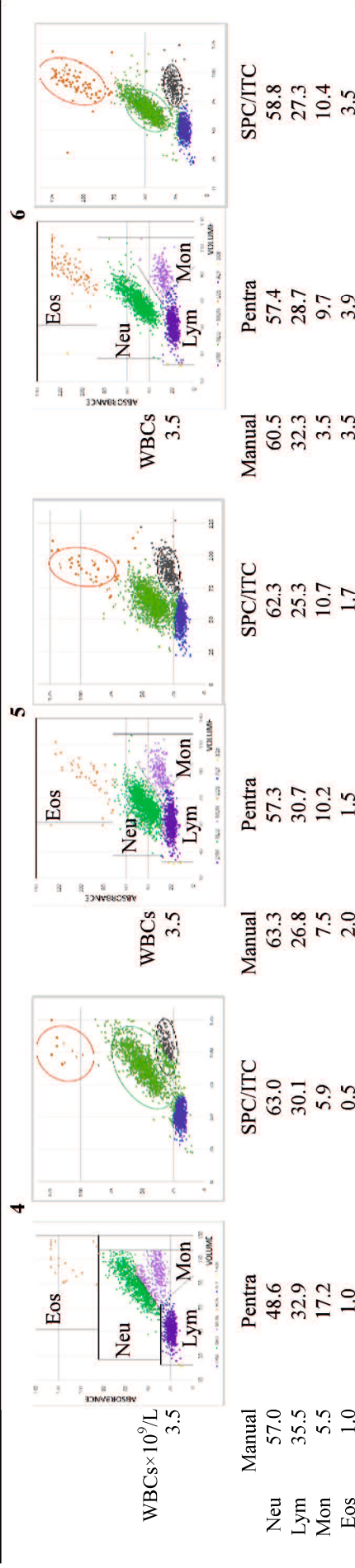
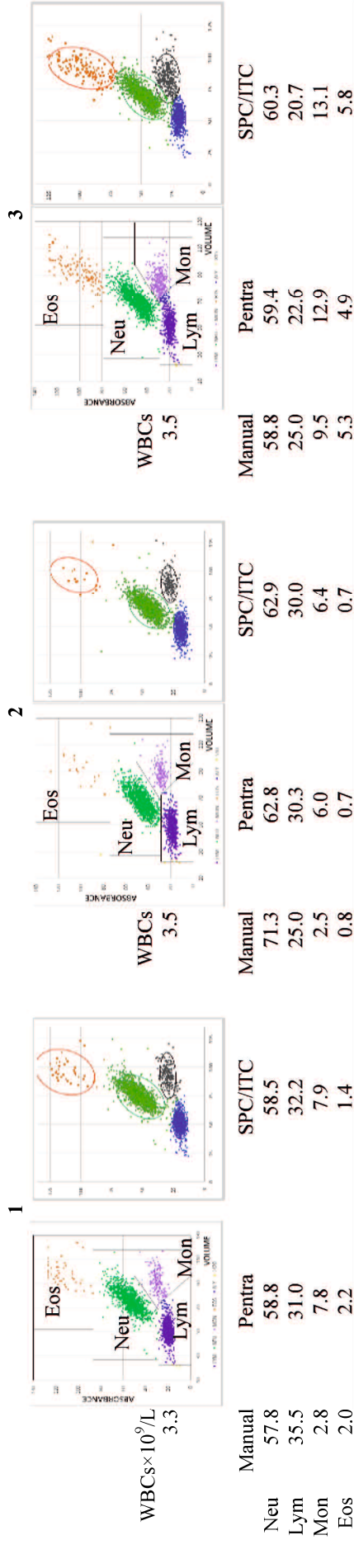
n は標本の大きさ、 s^2 はモデルの誤差項の分散推定量、 K はモデルに含まれる係数の数である。

第 1 項は、モデルの当てはまりを表すものですので、当てはまりが良いほど小さくなる。第 2 項は、「ペナルティー項」とも呼ばれ、パラメータ数が多いほど大きくなる。

- **目視法の結果と Pentra MS CPR の測定データおよび SPC/ITC 法適応図**

白血球分画について、健常者 118 名分および異常検体 109 名の合計 227 件分の分析結果を示す。左側の図は Pentra MS CPR で測定した固定した境界線によって計測したもので、右側の図は SPC/ITC 法によるものである。SPC/ITC 法の楕円は 95%信頼区間を示している。なお白血球の内、好塩基球に関しては Pentra MS CPR では別のチャンネルによって計測されるため、今回の計測には含まれていないことに注意。

健康者118名分の解析結果

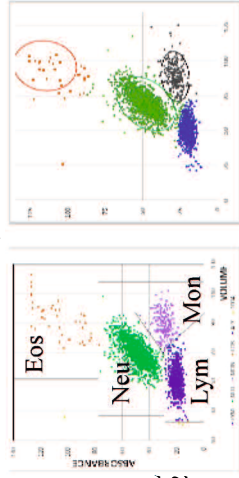


健常者118名分の解析結果

Case No.	WBCs × 10 ⁹ /L	Manual	Pentra	SPC/ITC	SPC/ITC	Manual	Pentra	SPC/ITC	SPC/ITC
10	3.8	49.0	51.2	54.8	54.8	48.0	46.4	48.8	48.8
11	3.8	41.5	37.8	36.8	37.6	44.3	41.2	39.5	39.5
12	3.8	2.8	6.1	4.9	6.8	5.3	10.8	10.6	10.6
13	3.9	5.3	4.6	3.4	0.5	1.5	1.4	1.2	1.2
14	3.9	Neu	52.0	54.9	58.1	63.0	58.4	61.7	68.1
15	3.9	Lym	38.5	32.7	29.5	25.8	31.3	28.8	17.9
16	4.0	Mon	5.5	9.7	10.5	7.3	11.2	10.5	10.5
17	4.0	Eos	2.8	2.4	1.8	3.3	3.8	3.5	3.5
18	4.1	Neu	61.8	53.6	56.3	62.3	59.5	63.5	63.5
19	4.1	Lym	33.3	37.7	36.6	32.5	32.8	29.4	29.4
20	4.1	Mon	3.8	6.9	6.2	2.3	5.0	5.0	5.0
21	4.1	Eos	1.0	1.4	0.8	1.8	2.3	2.1	2.1

健常者118名分の解析結果

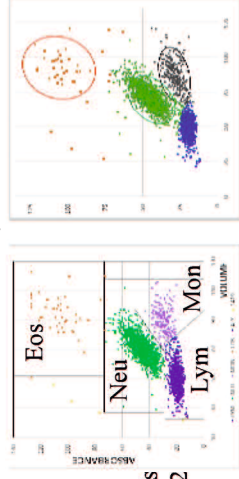
19



WBCs×10⁹/L
4.2

	Manual	SPC/ITC
Neu	56.8	57.3
Lym	34.8	33.9
Mon	6.0	7.9
Eos	2.0	1.0

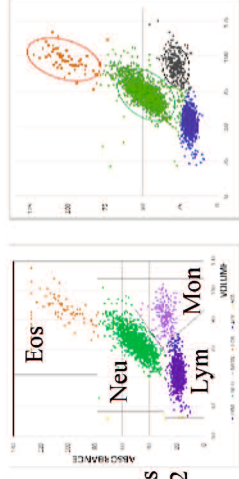
20



WBCs
4.2

	Manual	SPC/ITC
Neu	53.8	57.1
Lym	40.0	33.8
Mon	3.5	7.8
Eos	1.0	1.3

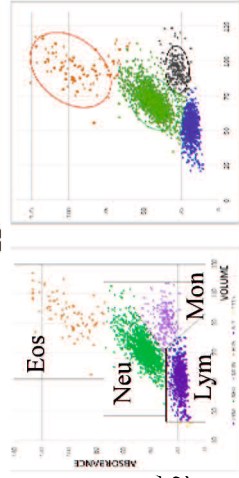
21



WBCs
4.2

	Manual	SPC/ITC
Neu	58.8	59.8
Lym	34.3	29.4
Mon	3.5	8.5
Eos	3.3	2.3

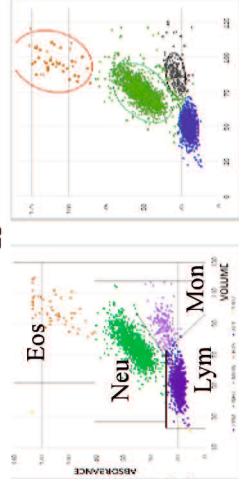
22



WBCs×10⁹/L
4.2

	Manual	SPC/ITC
Neu	55.3	54.4
Lym	37.8	35.9
Mon	2.5	6.6
Eos	2.3	3.0

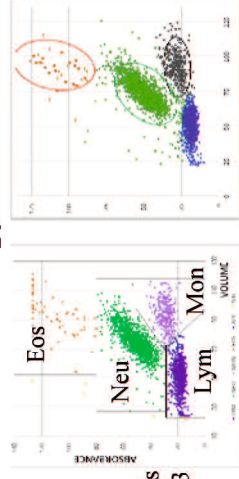
23



WBC
4.1

	Manual	SPC/ITC
Neu	45.5	50.3
Lym	46.5	42.7
Mon	5.3	5.6
Eos	1.8	1.4

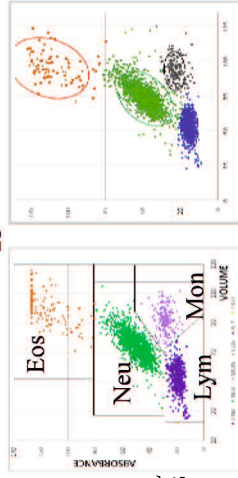
24



WBCs
4.3

	Manual	SPC/ITC
Neu	53.8	57.1
Lym	34.5	31.0
Mon	9.8	10.8
Eos	1.8	1.2

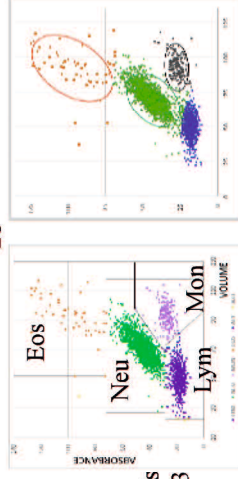
25



WBCs×10⁹/L
4.3

	Manual	SPC/ITC
Neu	62.3	61.8
Lym	31.0	29.8
Mon	2.0	5.7
Eos	3.5	2.8

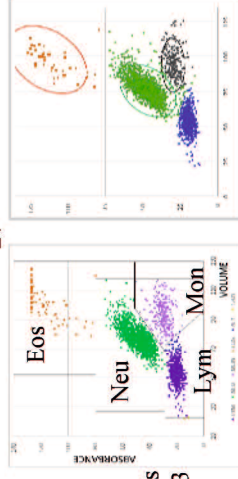
26



WBCs
4.3

	Manual	SPC/ITC
Neu	70.0	67.4
Lym	26.3	25.2
Mon	1.3	5.6
Eos	1.8	1.8

27

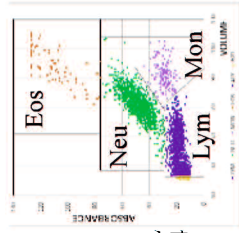


WBCs
4.3

	Manual	SPC/ITC
Neu	54.5	62.2
Lym	34.5	26.7
Mon	6.8	9.4
Eos	2.8	1.6

健常者118名分の解析結果

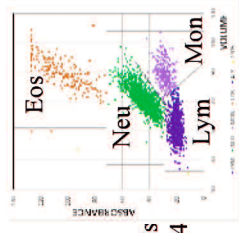
28



WBCs×10⁹/L
4.4

	Manual	Pentra	SPC/ITC
Neu	58.5	44.8	47.1
Lym	37.0	46.2	45.9
Mon	2.3	5.9	5.2
Eos	1.8	2.9	1.8

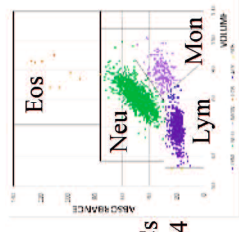
29



WBCs
4.4

	Manual	Pentra	SPC/ITC
Neu	55.5	48.7	52.0
Lym	32.8	34.1	30.4
Mon	3.8	11.4	11.8
Eos	6.5	5.7	5.8

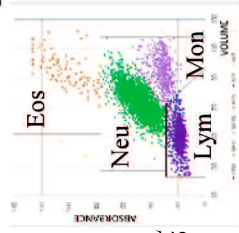
30



WBCs
4.4

	Manual	Pentra	SPC/ITC
Neu	69.8	68.0	69.8
Lym	25.0	24.4	22.6
Mon	3.8	7.2	7.3
Eos	0.8	0.3	0.3

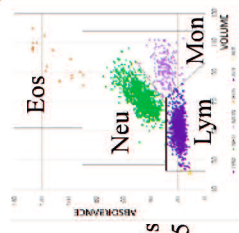
31



WBCs×10⁹/L
4.5

	Manual	Pentra	SPC/ITC
Neu	41.0	37.5	75.7
Lym	49.5	49.7	15.0
Mon	3.0	8.6	7.5
Eos	5.8	3.8	1.8

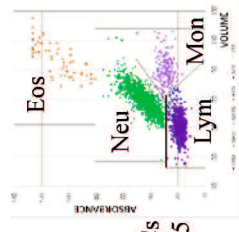
32



WBCs
4.5

	Manual	Pentra	SPC/ITC
Neu	50.8	49.7	53.8
Lym	44.5	42.5	40.1
Mon	3.3	6.8	5.7
Eos	1.3	0.7	0.4

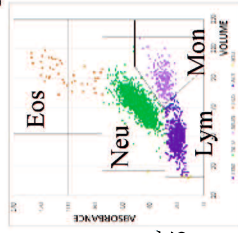
33



WBCs
4.5

	Manual	Pentra	SPC/ITC
Neu	56.0	54.6	58.9
Lym	39.0	35.7	32.6
Mon	3.3	6.8	6.8
Eos	1.5	2.5	1.7

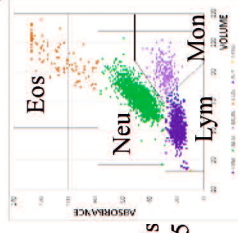
34



WBCs×10⁹/L
4.5

	Manual	Pentra	SPC/ITC
Neu	56.8	58.5	61.0
Lym	36.0	32.4	30.3
Mon	5.5	7.5	7.4
Eos	1.0	1.3	1.2

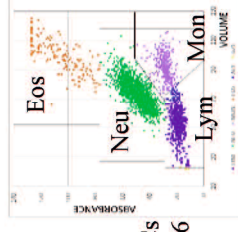
35



WBCs
4.5

	Manual	Pentra	SPC/ITC
Neu	55.0	57.1	59.7
Lym	37.8	32.8	31.8
Mon	2.8	6.6	5.9
Eos	3.0	3.2	2.6

36

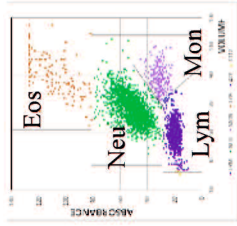


WBCs
4.6

	Manual	Pentra	SPC/ITC
Neu	67.8	58.7	62.3
Lym	24.8	28.2	25.3
Mon	2.0	8.9	8.6
Eos	4.3	4.2	3.8

健常者118名分の解析結果

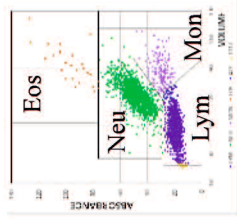
37



WBCs×10⁹/L
4.6

	Pentra	SPC/ITC
Manual	70.3	64.7
Neu	61.5	25.6
Lym	28.0	6.5
Mon	6.4	3.2
Eos	4.0	

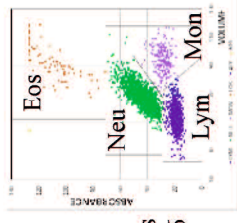
38



WBCs
4.6

	Pentra	SPC/ITC
Manual	60.3	56.8
Neu	55.6	35.9
Lym	38.2	6.4
Mon	5.1	0.9
Eos	0.9	

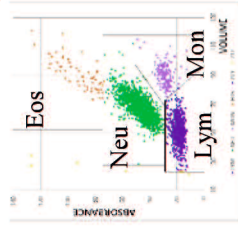
39



WBCs
4.6

	Manual	Pentra	SPC/ITC
Manual	67.3	61.4	66.0
Neu	25.8	28.6	25.0
Lym	4.0	7.2	7.2
Mon	1.0	2.5	1.8

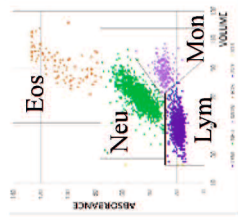
40



WBCs×10⁹/L
4.7

	Manual	Pentra	SPC/ITC
Manual	61.8	63.4	66.0
Neu	29.0	28.8	26.5
Lym	5.0	5.8	6.3
Mon	3.5	1.8	1.3

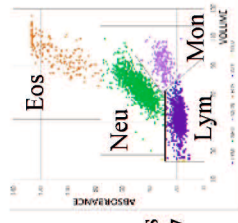
41



WBCs
4.7

	Manual	Pentra	SPC/ITC
Manual	52.3	56.2	55.6
Neu	39.0	33.8	29.5
Lym	5.8	7.0	11.4
Mon	1.8	2.9	3.5

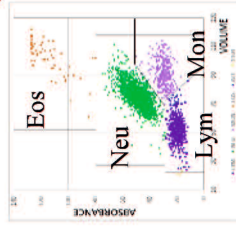
42



WBCs
4.7

	Manual	Pentra	SPC/ITC
Manual	46.3	50.2	53.6
Neu	41.0	36.9	35.5
Lym	3.3	7.7	6.9
Mon	5.5	5.0	4.0

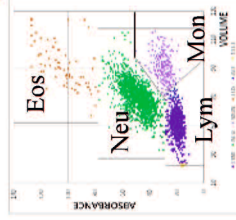
43



WBCs×10⁹/L
4.7

	Manual	Pentra	SPC/ITC
Manual	62.5	60.0	61.9
Neu	32.3	28.2	27.2
Lym	3.8	10.1	9.5
Mon	1.0	1.6	1.3

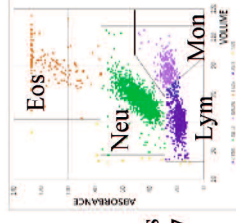
44



WBCs
4.7

	Manual	Pentra	SPC/ITC
Manual	57.8	52.9	54.9
Neu	34.8	38.0	36.3
Lym	4.3	6.6	6.8
Mon	2.3	2.4	2.0

45

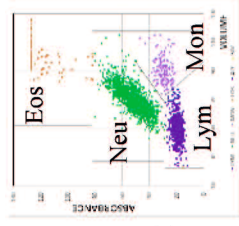


WBCs
4.7

	Manual	Pentra	SPC/ITC
Manual	59.0	59.0	62.3
Neu	30.5	27.5	24.9
Lym	5.8	8.4	10.0
Mon	3.8	4.7	2.8

健常者118名分の解析結果

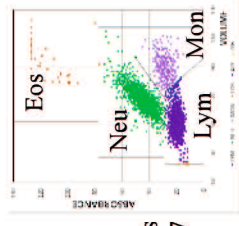
46



WBCs×10⁹/L
4.7

	Pentra	SPC/ITC
Manual	62.0	65.9
Neu	30.1	27.3
Lym	5.1	5.9
Mon	2.7	0.9
Eos		

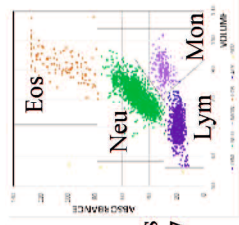
47



WBCs
4.7

	Pentra	SPC/ITC
Manual	45.5	49.4
Neu	43.3	38.7
Lym	9.7	11.4
Mon	1.4	0.6
Eos		

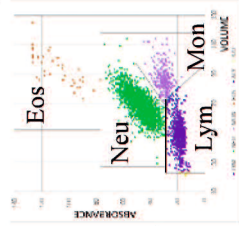
48



WBCs
4.7

	Manual	Pentra	SPC/ITC
Manual	62.5	63.9	66.0
Neu	28.8	25.5	24.9
Lym	4.3	6.8	6.2
Mon	3.8	3.7	2.9
Eos			

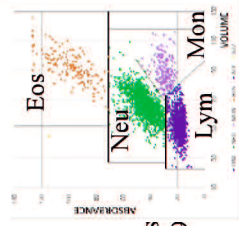
49



WBCs×10⁹/L
4.9

	Pentra	SPC/ITC
Manual	68.5	72.1
Neu	22.5	19.6
Lym	7.8	7.4
Mon	1.1	0.9
Eos		

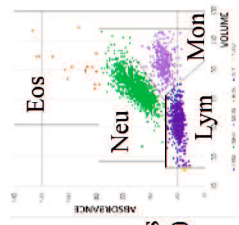
50



WBCs
4.9

	Manual	Pentra	SPC/ITC
Manual	55.8	49.8	52.5
Neu	37.3	39.3	35.6
Lym	3.3	6.2	7.5
Mon	2.8	4.5	4.4
Eos			

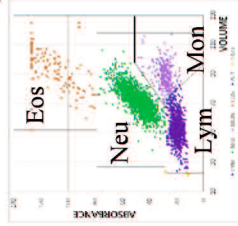
51



WBCs
5.0

	Manual	Pentra	SPC/ITC
Manual	60.5	52.9	55.1
Neu	31.0	33.3	31.8
Lym	5.8	13.2	12.6
Mon	0.8	0.5	0.4
Eos			

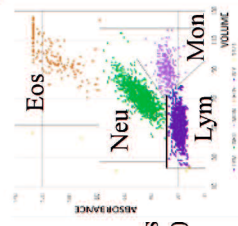
52



WBCs×10⁹/L
5.0

	Manual	Pentra	SPC/ITC
Manual	51.5	55.8	57.4
Neu	35.8	32.2	30.8
Lym	6.8	8.2	8.9
Mon	3.5	3.5	3.0
Eos			

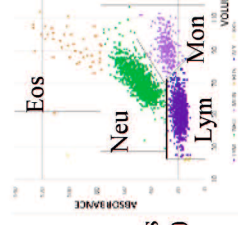
53



WBCs
5.0

	Manual	Pentra	SPC/ITC
Manual	47.3	41.7	44.0
Neu	47.5	43.9	43.2
Lym	1.8	9.4	9.8
Mon	3.0	4.9	3.0
Eos			

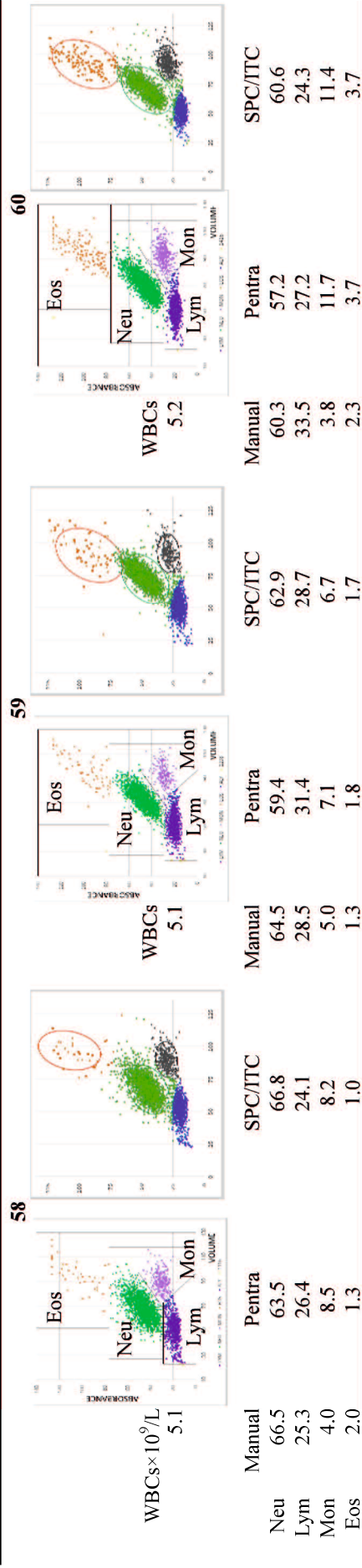
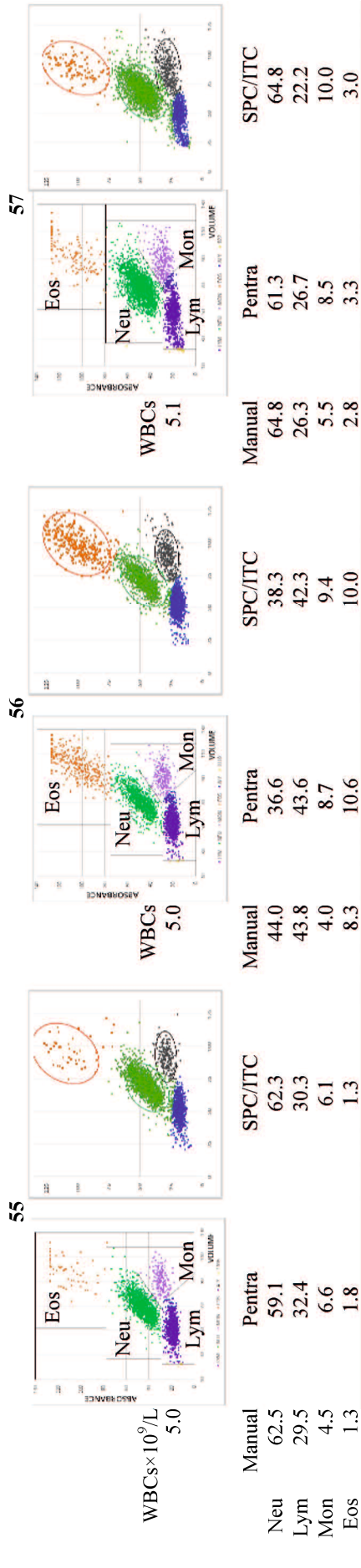
54



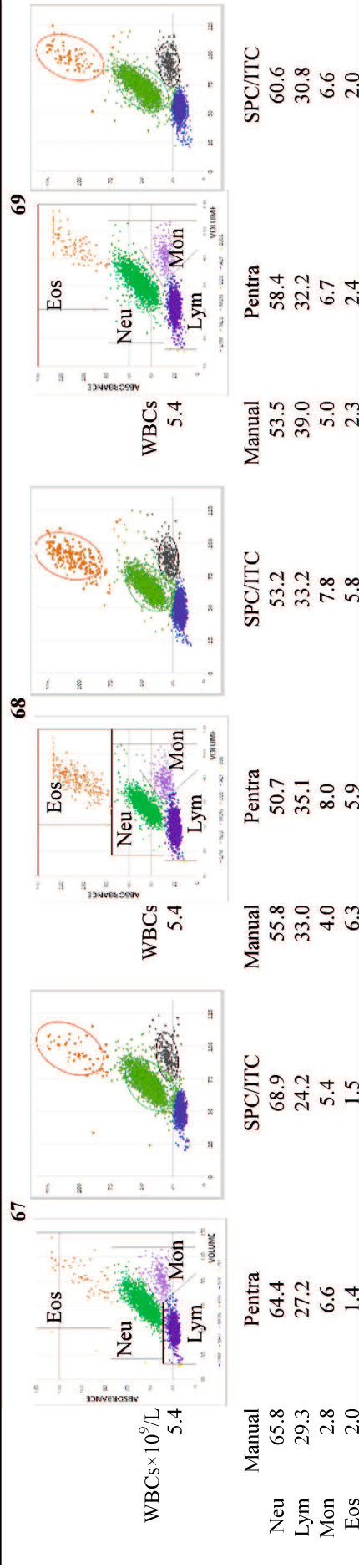
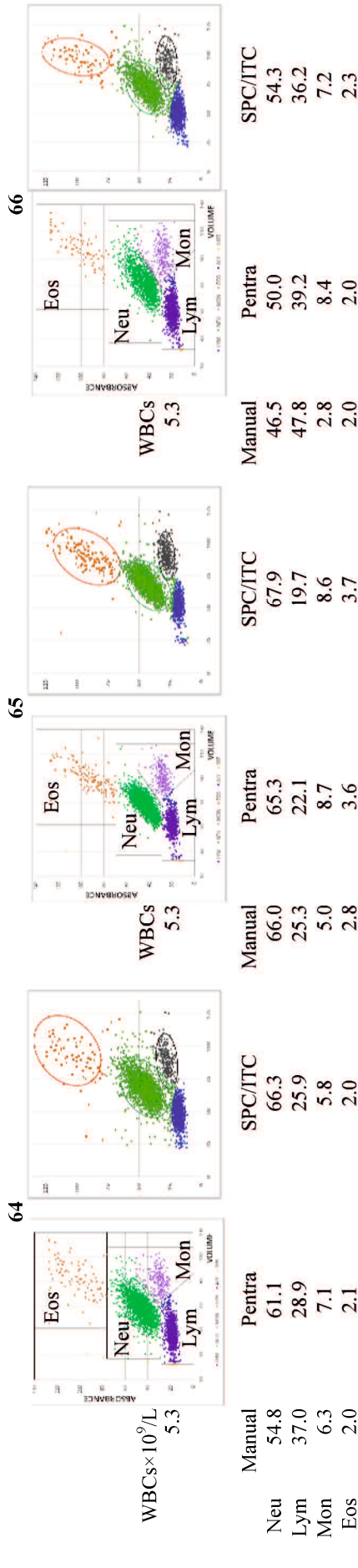
WBCs
5.0

	Manual	Pentra	SPC/ITC
Manual	48.0	46.7	48.5
Neu	45.0	42.9	41.3
Lym	4.0	9.1	9.0
Mon	0.8	1.1	1.1
Eos			

健康者118名分の解析結果



健康者118名分の解析結果

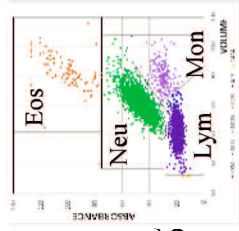


健康者118名分の解析結果

73		WBCs×10 ⁹ /L 5.5	74		WBCs 5.6	75		WBCs 5.7
Manual	Pentra	SPC/ITC	Manual	Pentra	SPC/ITC	Manual	Pentra	SPC/ITC
Neu	49.0	53.0	67.5	60.3	62.6	65.5	65.3	67.7
Lym	42.4	39.2	27.0	27.3	26.0	28.0	25.8	24.1
Mon	6.3	6.0	3.0	9.5	9.8	3.0	6.9	6.4
Eos	2.0	1.8	2.0	2.7	1.6	2.0	1.8	1.8
76		WBCs×10 ⁹ /L 5.7	77		WBCs 5.7	78		WBCs 5.7
Manual	Pentra	SPC/ITC	Manual	Pentra	SPC/ITC	Manual	Pentra	SPC/ITC
Neu	71.3	75.4	58.0	57.7	61.3	64.3	63.5	67.2
Lym	22.0	20.0	35.3	34.7	32.1	25.8	24.7	22.0
Mon	4.9	3.1	4.5	6.3	5.9	6.3	7.3	7.9
Eos	1.7	1.5	1.0	1.0	0.7	3.0	4.3	2.9
79		WBCs×10 ⁹ /L 5.8	80		WBCs 5.8	81		WBCs 5.8
Manual	Pentra	SPC/ITC	Manual	Pentra	SPC/ITC	Manual	Pentra	SPC/ITC
Neu	52.4	55.0	49.5	42.5	44.9	54.0	54.6	57.1
Lym	37.3	36.0	39.5	40.7	38.2	40.0	36.7	34.2
Mon	7.4	7.3	5.3	11.9	13.2	3.8	7.3	7.6
Eos	2.8	1.7	5.0	4.6	3.7	1.5	1.3	1.0

健康者118名分の解析結果

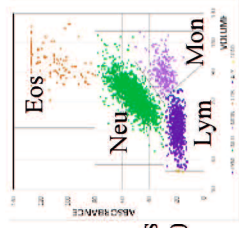
82



WBCs×10⁹/L
5.9

	Manual	Pentra	SPC/ITC
Neu	70.8	66.5	69.6
Lym	26.8	25.2	23.6
Mon	1.5	5.9	4.8
Eos	1.0	2.2	2.0

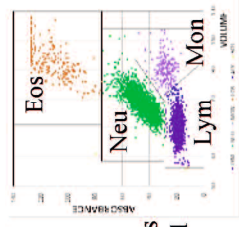
83



WBCs
6.0

	Manual	Pentra	SPC/ITC
Neu	69.5	66.3	68.3
Lym	23.0	23.4	22.4
Mon	5.0	7.5	7.5
Eos	1.3	2.5	1.8

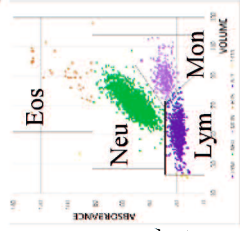
84



WBCs
6.1

	Manual	Pentra	SPC/ITC
Neu	69.5	60.6	65.1
Lym	23.8	27.9	24.5
Mon	1.8	6.0	6.3
Eos	3.8	5.4	4.0

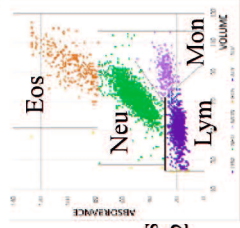
85



WBCs×10⁹/L
6.1

	Manual	Pentra	SPC/ITC
Neu	65.0	65.1	66.5
Lym	27.8	25.6	24.4
Mon	4.8	8.2	8.1
Eos	1.3	0.8	0.9

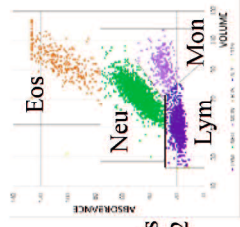
86



WBCs
6.2

	Manual	Pentra	SPC/ITC
Neu	60.3	58.4	59.6
Lym	26.8	28.1	26.6
Mon	5.3	6.5	6.6
Eos	6.0	6.2	7.2

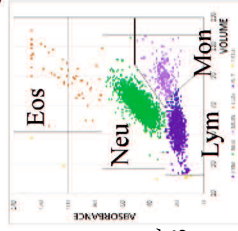
87



WBCs
6.2

	Manual	Pentra	SPC/ITC
Neu	56.8	56.9	61.4
Lym	30.8	28.4	25.8
Mon	6.5	7.7	7.3
Eos	4.8	6.5	5.5

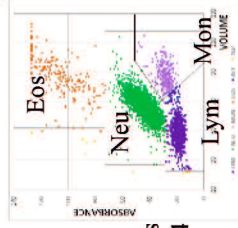
88



WBCs×10⁹/L
6.3

	Manual	Pentra	SPC/ITC
Neu	53.3	52.3	54.6
Lym	44.3	40.0	37.7
Mon	1.8	6.0	6.7
Eos	0.8	1.2	1.0

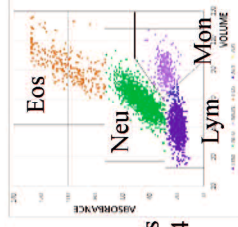
89



WBCs
6.4

	Manual	Pentra	SPC/ITC
Neu	65.0	58.1	61.5
Lym	27.5	29.5	27.1
Mon	4.0	7.4	7.4
Eos	3.0	4.6	3.9

90

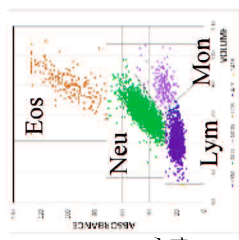


WBCs
6.4

	Manual	Pentra	SPC/ITC
Neu	51.5	50.7	54.1
Lym	36.8	36.0	32.7
Mon	6.3	7.3	7.5
Eos	4.8	5.7	5.8

健康者118名分の解析結果

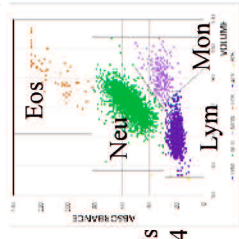
91



WBCs×10⁹/L
6.4

	Manual	Pentra	SPC/ITC
Neu	57.5	57.5	61.6
Lym	36.3	32.3	28.7
Mon	2.5	5.9	5.5
Eos	2.5	4.1	4.1

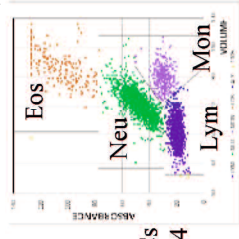
92



WBCs
6.4

	Manual	Pentra	SPC/ITC
Neu	65.0	63.0	64.4
Lym	28.3	29.2	28.5
Mon	3.0	6.0	6.0
Eos	2.0	1.5	1.0

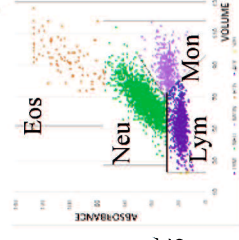
93



WBCs
6.4

	Manual	Pentra	SPC/ITC
Neu	48.8	45.2	49.4
Lym	39.8	40.3	37.6
Mon	5.3	9.4	9.1
Eos	5.0	4.6	3.9

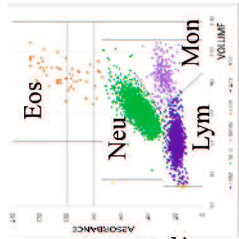
94



WBCs×10⁹/L
6.5

	Manual	Pentra	SPC/ITC
Neu	61.5	57.7	63.0
Lym	29.0	27.4	23.2
Mon	7.0	12.5	12.0
Eos	2.3	2.1	1.8

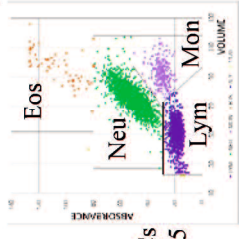
95



WBC
6.5

	Manual	Pentra	SPC/ITC
Neu	57.5	57.7	60.4
Lym	37.5	33.0	31.3
Mon	4.0	7.4	7.0
Eos	0.8	1.4	1.3

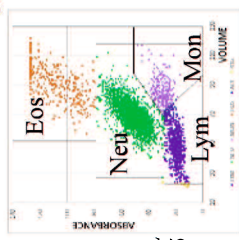
96



WBCs
6.5

	Manual	Pentra	SPC/ITC
Neu	59.5	59.8	62.7
Lym	35.8	32.4	30.9
Mon	2.5	6.4	5.5
Eos	1.0	1.2	1.0

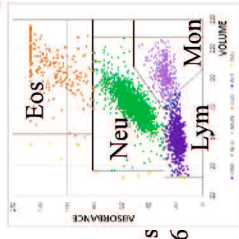
97



WBCs×10⁹/L
6.5

	Manual	Pentra	SPC/ITC
Neu	67.5	63.8	65.8
Lym	24.0	22.6	20.5
Mon	3.0	6.7	7.7
Eos	4.5	6.7	5.9

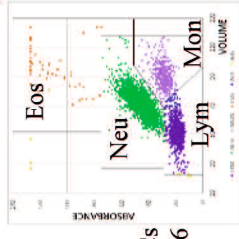
98



WBCs
6.6

	Manual	Pentra	SPC/ITC
Neu	54.0	55.7	57.8
Lym	35.0	30.3	29.4
Mon	5.8	8.8	8.6
Eos	4.0	4.7	4.2

99



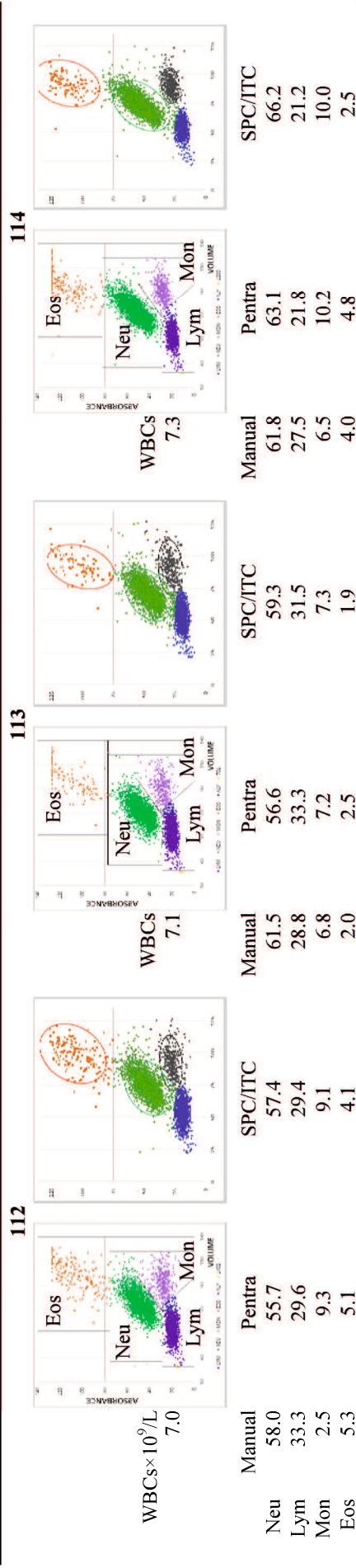
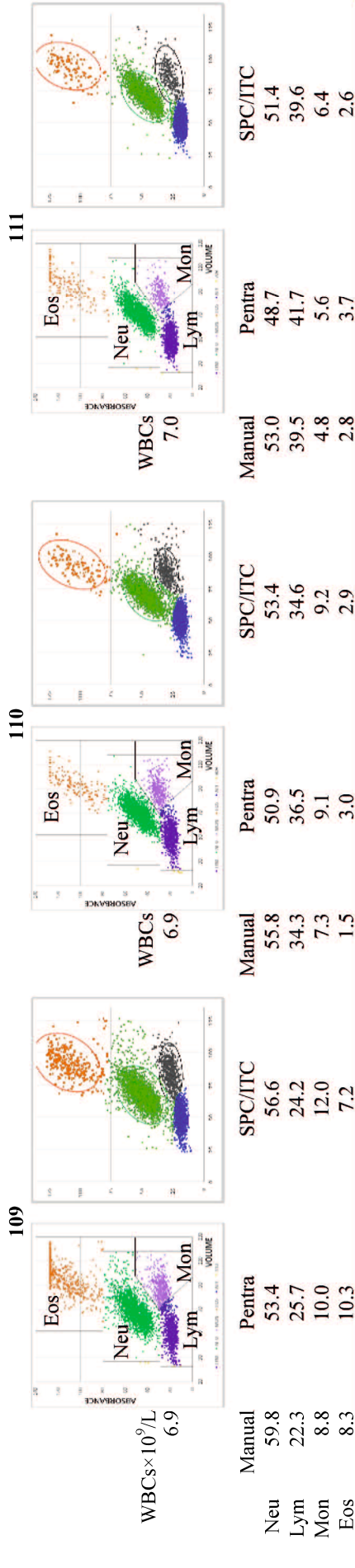
WBCs
6.6

	Manual	Pentra	SPC/ITC
Neu	64.8	65.1	68.5
Lym	25.8	24.0	21.5
Mon	7.0	9.4	9.4
Eos	1.0	1.2	0.7

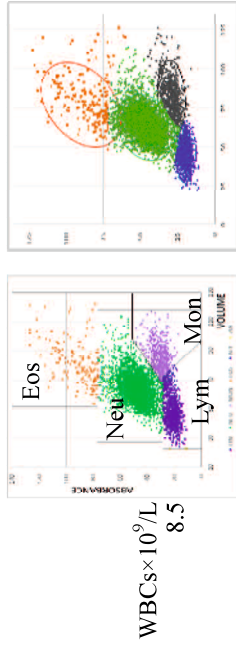
健康者118名分の解析結果

100		101		102	
WBCs×10 ⁹ /L 6.6	WBCs 6.6	WBCs×10 ⁹ /L 6.6	WBCs 6.6	WBCs 6.7	WBCs×10 ⁹ /L 6.6
Manual	Manual	Manual	Manual	Manual	Manual
Neu 53.5	Neu 57.3	Neu 58.1	Neu 60.6	Neu 64.4	Neu 66.5
Lym 38.0	Lym 35.8	Lym 32.2	Lym 30.9	Lym 27.5	Lym 25.6
Mon 3.0	Mon 3.3	Mon 6.1	Mon 5.6	Mon 5.5	Mon 5.3
Eos 4.5	Eos 3.0	Eos 3.4	Eos 2.9	Eos 2.4	Eos 2.6
Pentra 56.8	Pentra 57.3	Pentra 58.1	Pentra 60.6	Pentra 64.4	Pentra 66.5
SPC/ITC 27.2	SPC/ITC 27.2	SPC/ITC 32.2	SPC/ITC 30.9	SPC/ITC 27.5	SPC/ITC 25.6
SPC/ITC 9.2	SPC/ITC 9.2	SPC/ITC 9.2	SPC/ITC 5.6	SPC/ITC 5.5	SPC/ITC 5.3
SPC/ITC 3.7	SPC/ITC 3.7	SPC/ITC 3.4	SPC/ITC 2.9	SPC/ITC 2.4	SPC/ITC 2.6
103		104		105	
WBCs×10 ⁹ /L 6.7	WBCs 6.8	WBCs×10 ⁹ /L 6.7	WBCs 6.8	WBCs 6.8	WBCs×10 ⁹ /L 6.7
Manual	Manual	Manual	Manual	Manual	Manual
Neu 57.3	Neu 71.8	Neu 68.7	Neu 71.6	Neu 59.8	Neu 62.5
Lym 37.5	Lym 22.0	Lym 24.5	Lym 22.4	Lym 32.1	Lym 30.3
Mon 3.0	Mon 4.8	Mon 5.3	Mon 5.3	Mon 7.3	Mon 6.7
Eos 1.8	Eos 1.0	Eos 1.2	Eos 0.7	Eos 0.6	Eos 0.4
Pentra 57.3	Pentra 71.8	Pentra 68.7	Pentra 71.6	Pentra 59.8	Pentra 62.5
SPC/ITC 33.4	SPC/ITC 33.4	SPC/ITC 24.5	SPC/ITC 22.4	SPC/ITC 32.1	SPC/ITC 30.3
SPC/ITC 5.7	SPC/ITC 5.7	SPC/ITC 5.3	SPC/ITC 5.3	SPC/ITC 7.3	SPC/ITC 6.7
SPC/ITC 1.4	SPC/ITC 1.4	SPC/ITC 1.2	SPC/ITC 0.7	SPC/ITC 0.6	SPC/ITC 0.4
106		107		108	
WBCs×10 ⁹ /L 6.8	WBCs 6.8	WBCs×10 ⁹ /L 6.8	WBCs 6.8	WBCs 6.9	WBCs×10 ⁹ /L 6.8
Manual	Manual	Manual	Manual	Manual	Manual
Neu 66.0	Neu 53.5	Neu 56.0	Neu 59.3	Neu 54.4	Neu 56.3
Lym 26.0	Lym 32.8	Lym 31.5	Lym 29.1	Lym 33.5	Lym 33.0
Mon 3.8	Mon 3.8	Mon 5.0	Mon 6.1	Mon 6.2	Mon 6.1
Eos 3.5	Eos 7.3	Eos 7.1	Eos 5.5	Eos 5.6	Eos 4.6
Pentra 57.5	Pentra 53.5	Pentra 56.0	Pentra 59.3	Pentra 54.4	Pentra 56.3
SPC/ITC 32.0	SPC/ITC 32.8	SPC/ITC 31.5	SPC/ITC 29.1	SPC/ITC 33.5	SPC/ITC 33.0
SPC/ITC 7.7	SPC/ITC 3.8	SPC/ITC 5.0	SPC/ITC 6.1	SPC/ITC 6.2	SPC/ITC 6.1
SPC/ITC 2.7	SPC/ITC 7.3	SPC/ITC 7.1	SPC/ITC 5.5	SPC/ITC 5.6	SPC/ITC 4.6

健常者118名分の解析結果

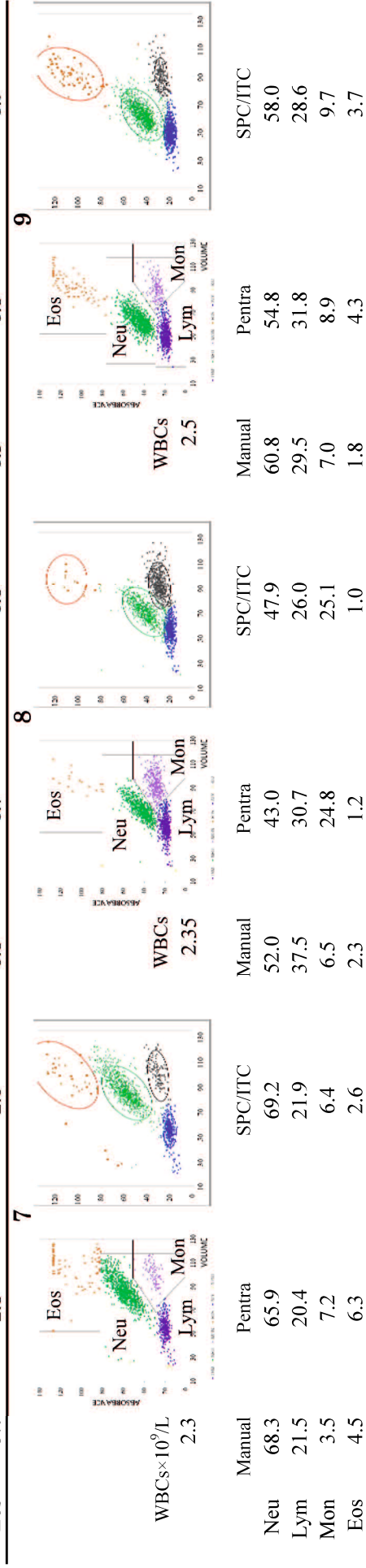
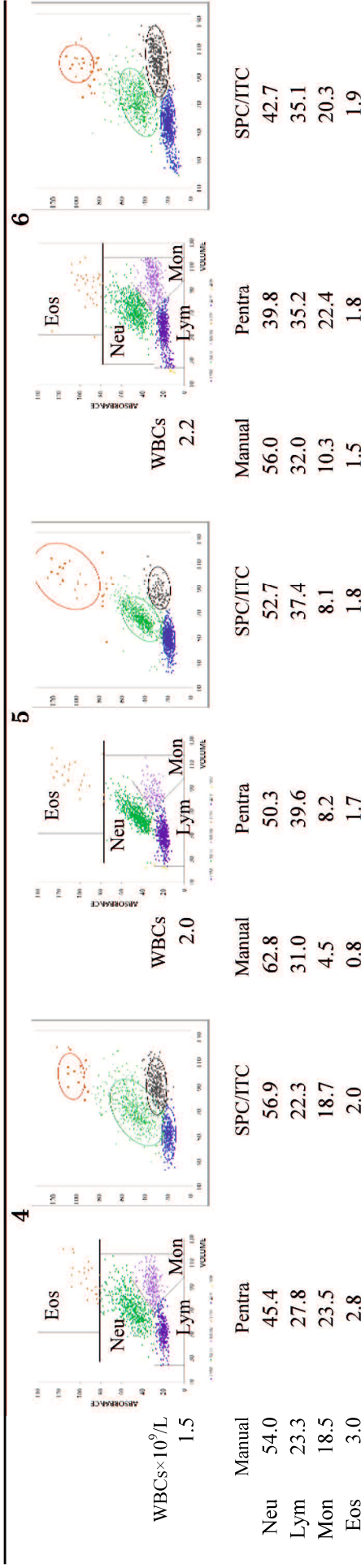
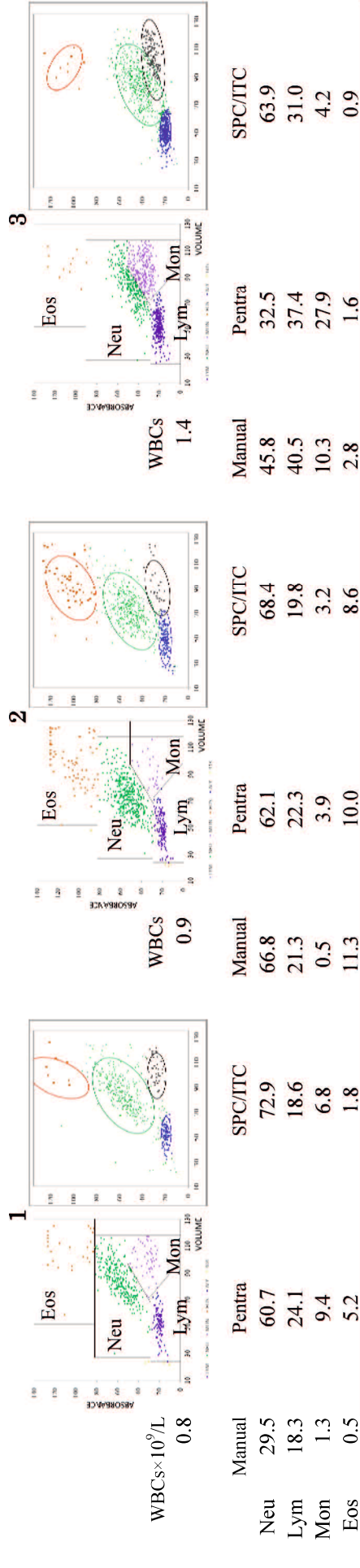


118

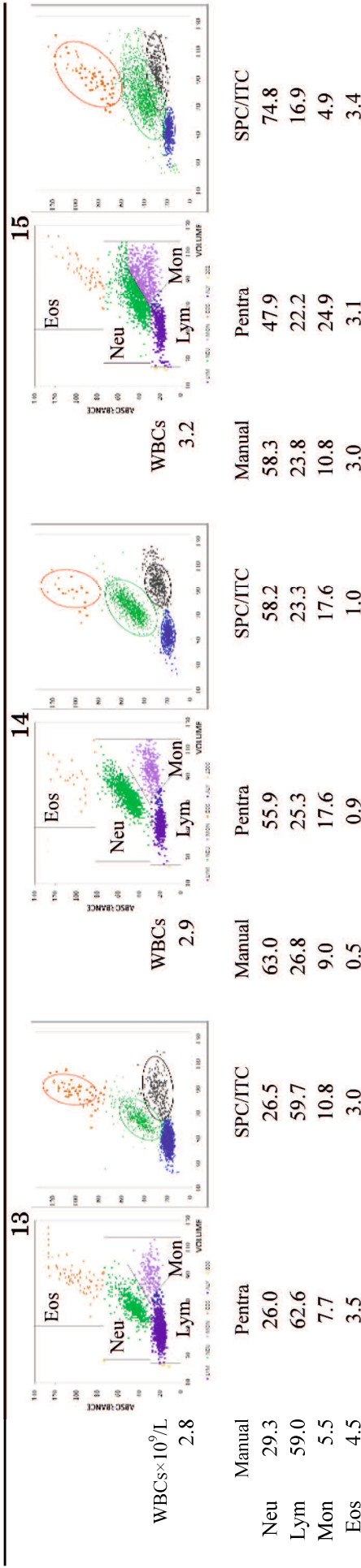
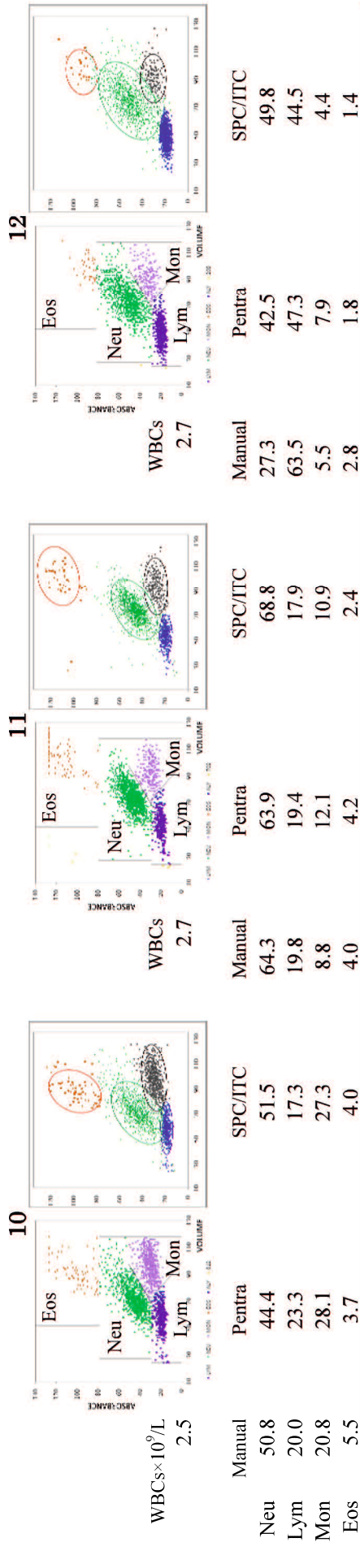


	Manual	Pentra	SPC/ITC
Neu	64.3	66.7	68.0
Lym	21.0	21.2	17.1
Mon	10.5	8.8	11.1
Eos	2.3	2.8	3.8

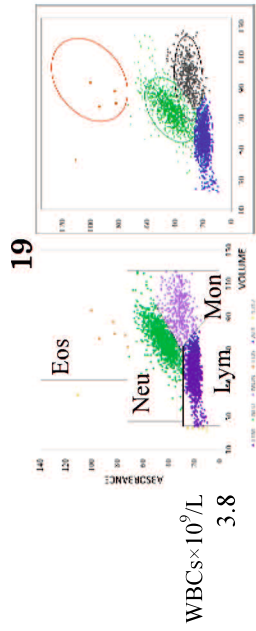
異常検体109名分の解析結果



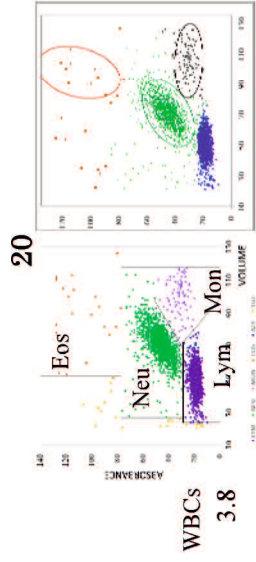
異常検体109名分の解析結果



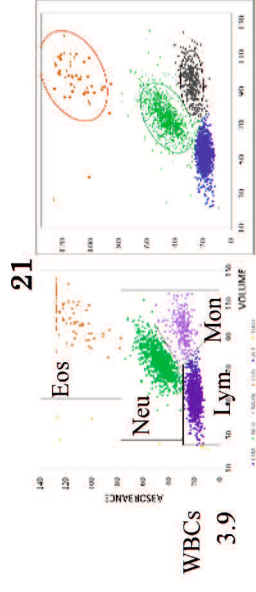
異常検体109名分の解析結果



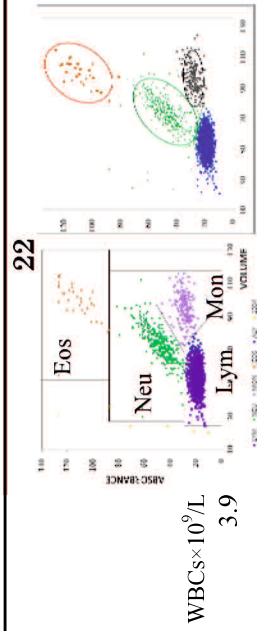
	Manual	SPC/ITC
Neu	62.0	71.0
Lym	30.5	9.2
Mon	5.8	18.9
Eos	0.5	0.9



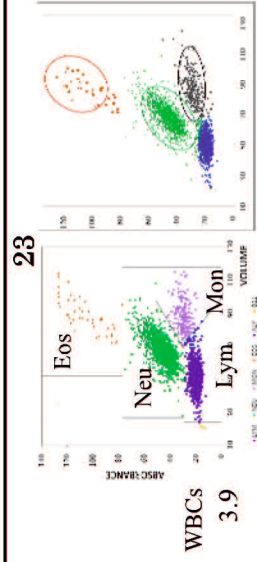
	Manual	SPC/ITC
Neu	43.8	45.5
Lym	52.5	49.4
Mon	3.0	3.2
Eos	0.5	1.0



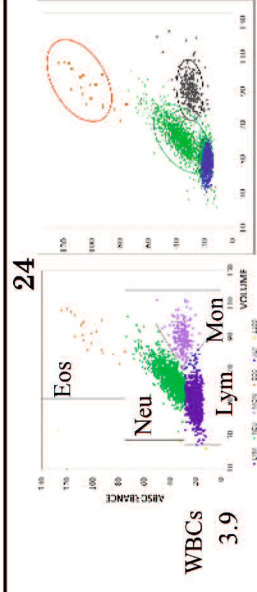
	Manual	SPC/ITC
Neu	37.0	39.5
Lym	54.3	47.0
Mon	4.0	10.1
Eos	3.3	3.2



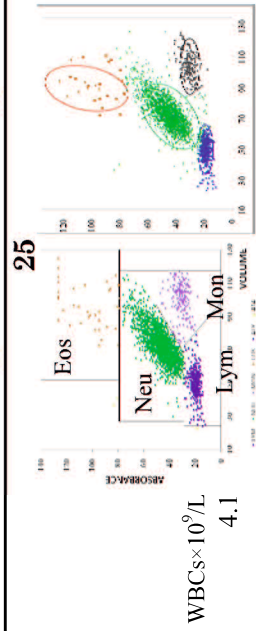
	Manual	SPC/ITC
Neu	12.5	14.4
Lym	81.8	77.5
Mon	3.0	7.1
Eos	2.0	1.0



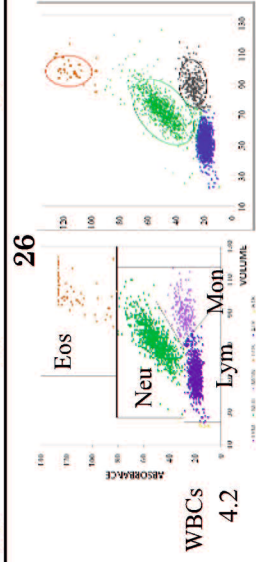
	Manual	SPC/ITC
Neu	56.8	53.3
Lym	36.3	37.6
Mon	5.0	7.3
Eos	0.5	1.6



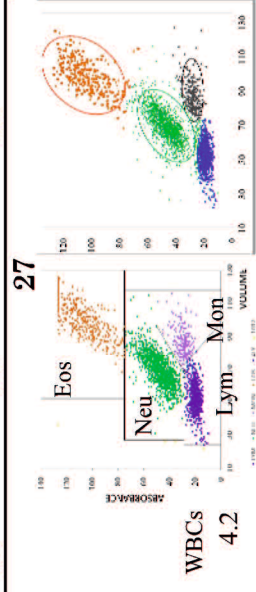
	Manual	SPC/ITC
Neu	51.5	45.7
Lym	41.3	44.6
Mon	6.0	8.6
Eos	0.5	0.8



	Manual	SPC/ITC
Neu	80.5	81.8
Lym	14.8	11.4
Mon	3.5	5.7
Eos	0.8	1.1

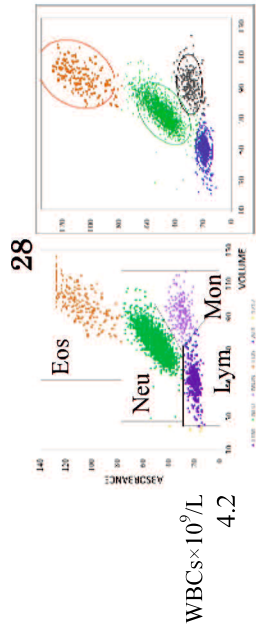


	Manual	SPC/ITC
Neu	43.3	40.0
Lym	50.5	49.4
Mon	3.3	7.7
Eos	2.0	2.6

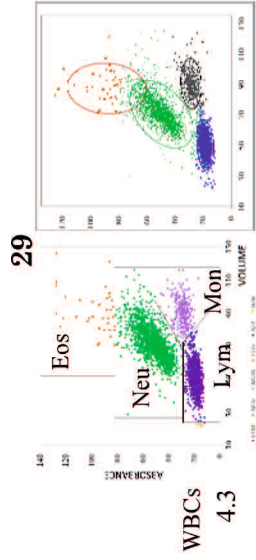


	Manual	SPC/ITC
Neu	49.8	51.1
Lym	33.8	29.9
Mon	5.8	7.5
Eos	9.8	11.2

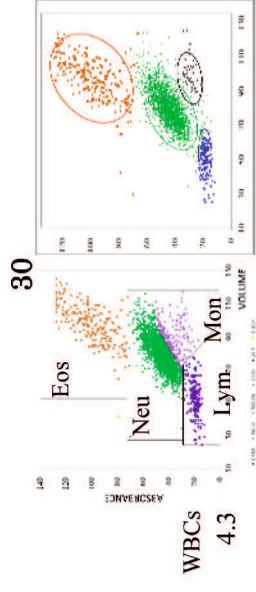
異常検体109名分の解析結果



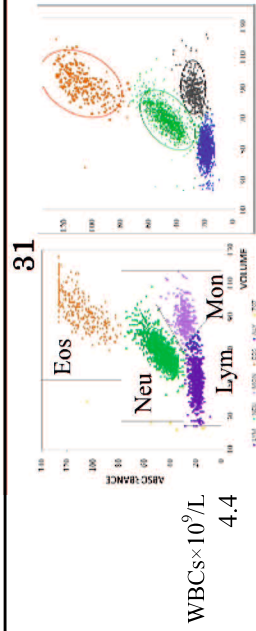
	Manual	Pentra	SPC/ITC
Neu	69.3	67.4	70.3
Lym	16.3	18.0	15.8
Mon	7.3	7.4	7.6
Eos	6.0	7.2	6.3



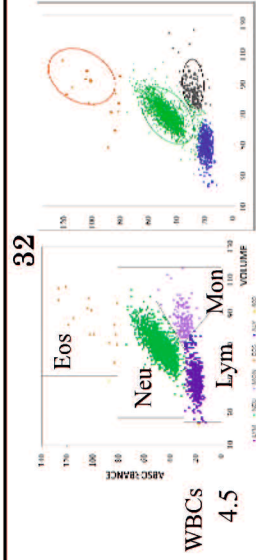
	Manual	Pentra	SPC/ITC
Neu	54.5	51.6	54.1
Lym	38.0	38.1	35.9
Mon	5.3	9.2	8.5
Eos	0.5	0.9	1.4



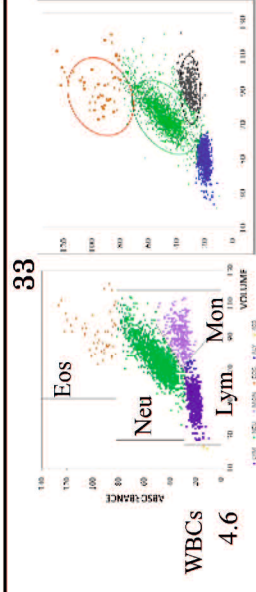
	Manual	Pentra	SPC/ITC
Neu	82.8	74.2	84.8
Lym	9.5	7.8	6.0
Mon	0.5	9.6	1.2
Eos	7.0	8.0	8.0



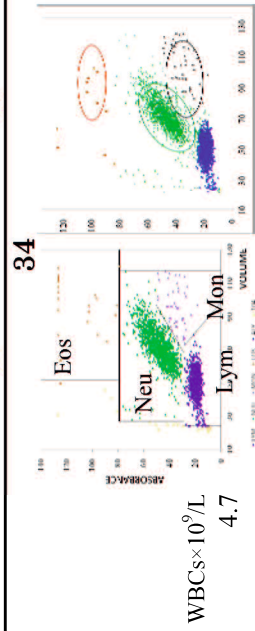
	Manual	Pentra	SPC/ITC
Neu	35.3	40.0	42.8
Lym	45.0	39.2	37.4
Mon	8.0	10.9	11.5
Eos	9.8	9.6	8.3



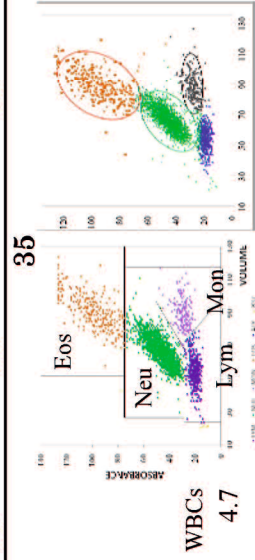
	Manual	Pentra	SPC/ITC
Neu	79.0	69.3	74.1
Lym	15.8	22.3	20.0
Mon	4.5	7.7	5.4
Eos	0.5	0.6	0.5



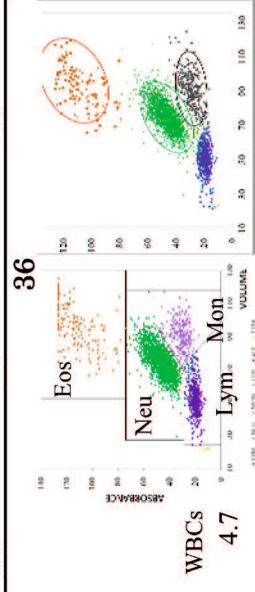
	Manual	Pentra	SPC/ITC
Neu	74.0	60.6	63.9
Lym	19.8	26.9	23.6
Mon	5.5	11.0	8.9
Eos	0.8	1.5	3.6



	Manual	Pentra	SPC/ITC
Neu	59.8	55.1	57.7
Lym	39.0	42.0	40.6
Mon	1.0	1.6	1.5
Eos	0.3	0.8	0.2

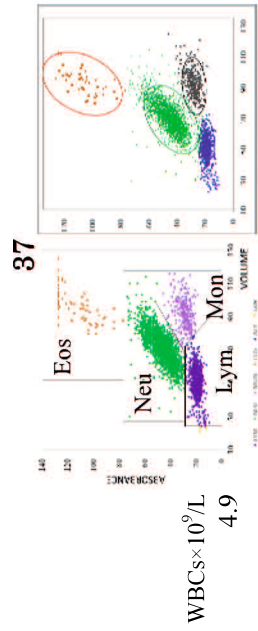


	Manual	Pentra	SPC/ITC
Neu	75.8	70.9	74.6
Lym	9.5	15.3	11.1
Mon	3.8	5.8	6.2
Eos	8.3	7.7	8.1

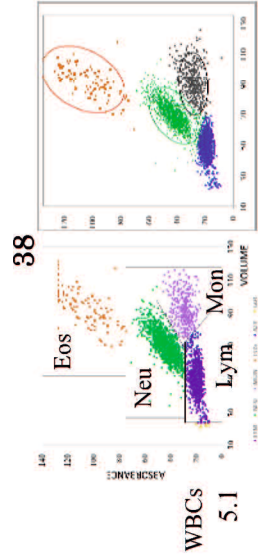


	Manual	Pentra	SPC/ITC
Neu	73.3	68.0	73.0
Lym	18.0	16.9	15.5
Mon	4.0	9.1	7.8
Eos	4.3	6.0	3.6

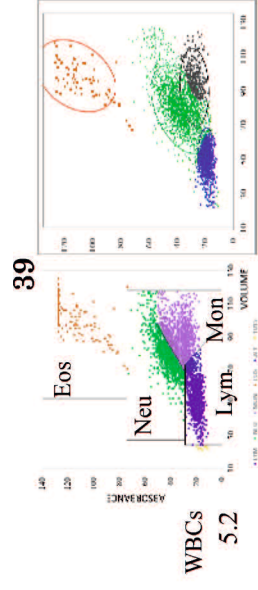
異常検体109名分の解析結果



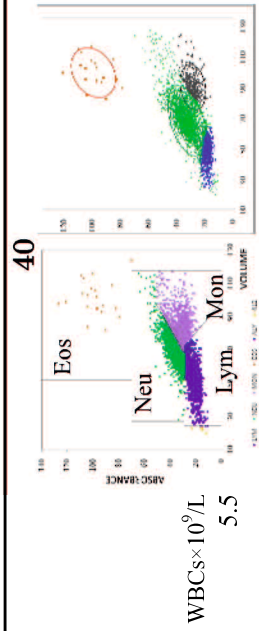
	Manual	SPC/ITC
Neu	73.0	73.0
Lym	18.8	16.0
Mon	5.3	9.5
Eos	1.8	1.5



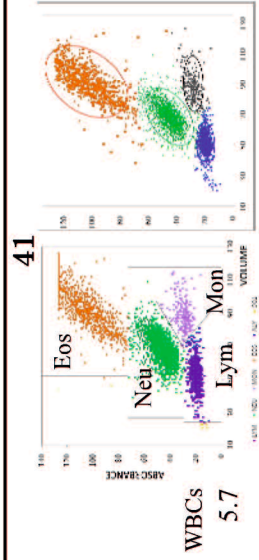
	Manual	SPC/ITC
Neu	35.8	42.8
Lym	51.3	42.0
Mon	9.3	11.5
Eos	2.3	3.3



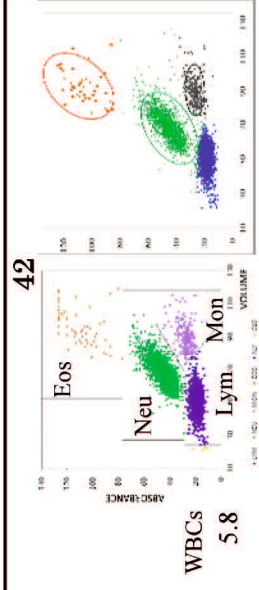
	Manual	SPC/ITC
Neu	40.8	30.5
Lym	33.3	34.6
Mon	14.0	30.5
Eos	4.0	3.6



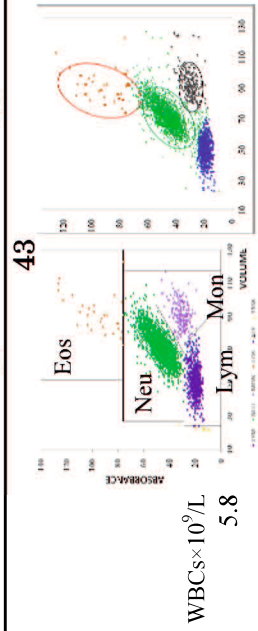
	Manual	SPC/ITC
Neu	68.3	77.3
Lym	20.3	16.8
Mon	7.3	5.3
Eos	0.8	0.5



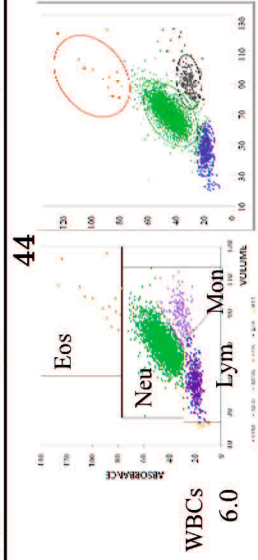
	Manual	SPC/ITC
Neu	52.5	48.9
Lym	16.3	19.0
Mon	2.3	5.1
Eos	28.3	26.8



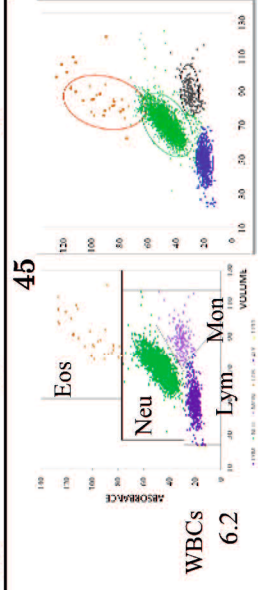
	Manual	SPC/ITC
Neu	61.0	51.5
Lym	35.3	40.4
Mon	1.8	6.3
Eos	0.5	1.4



	Manual	SPC/ITC
Neu	74.3	73.9
Lym	20.0	19.4
Mon	3.0	5.7
Eos	0.8	1.0

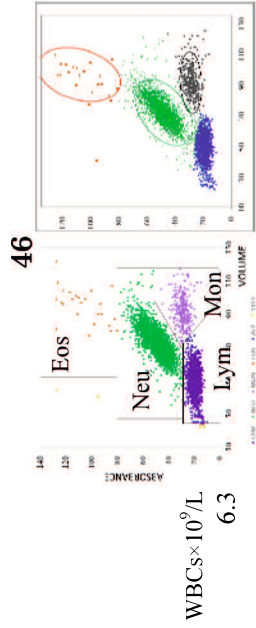


	Manual	SPC/ITC
Neu	88.0	83.3
Lym	8.0	10.9
Mon	2.8	5.3
Eos	1.0	0.4

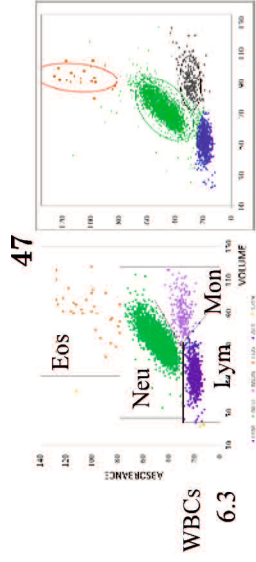


	Manual	SPC/ITC
Neu	82.8	74.5
Lym	13.0	19.1
Mon	2.3	5.3
Eos	1.5	1.1

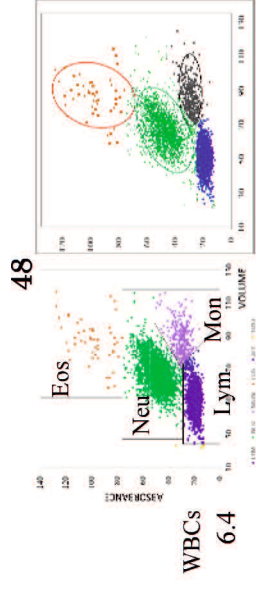
異常検体109名分の解析結果



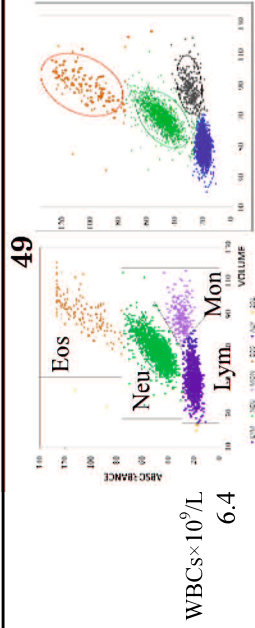
WBCs×10 ⁹ /L		SPC/ITC	
Manual	Pentra	Manual	SPC/ITC
Neu	58.8	61.4	
Lym	20.8	30.3	
Mon	1.8	7.7	
Eos	0.5	0.5	



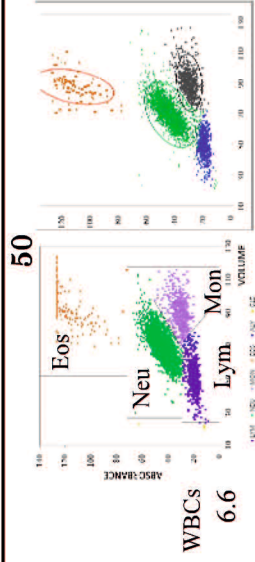
WBCs		SPC/ITC	
Manual	Pentra	Manual	SPC/ITC
Neu	73.8	77.3	
Lym	20.2	17.6	
Mon	5.1	4.6	
Eos	0.8	0.6	



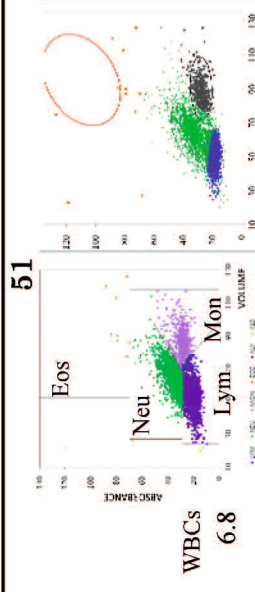
WBCs		SPC/ITC	
Manual	Pentra	Manual	SPC/ITC
Neu	50.4	55.4	
Lym	40.1	37.6	
Mon	8.2	5.9	
Eos	1.0	1.0	



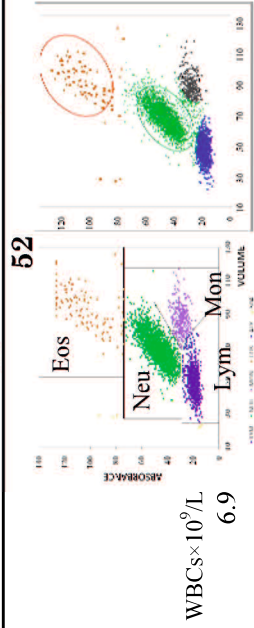
WBCs×10 ⁹ /L		SPC/ITC	
Manual	Pentra	Manual	SPC/ITC
Neu	36.0	43.5	
Lym	52.5	46.4	
Mon	5.3	6.7	
Eos	3.8	3.3	



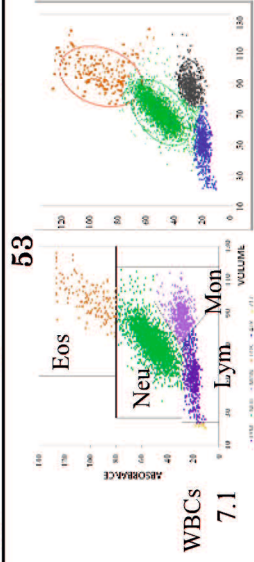
WBCs		SPC/ITC	
Manual	Pentra	Manual	SPC/ITC
Neu	67.0	72.8	
Lym	15.2	12.2	
Mon	13.5	13.1	
Eos	3.9	1.8	



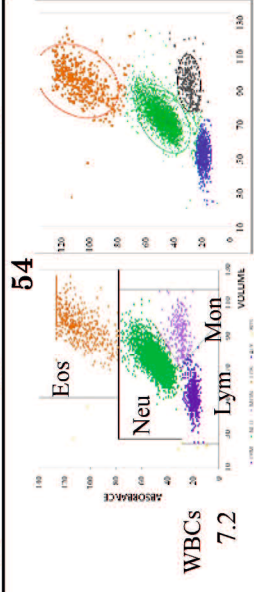
WBCs		SPC/ITC	
Manual	Pentra	Manual	SPC/ITC
Neu	31.2	50.6	
Lym	55.0	40.1	
Mon	13.1	9.3	
Eos	1.0	0.0	



WBCs×10 ⁹ /L		SPC/ITC	
Manual	Pentra	Manual	SPC/ITC
Neu	76.3	69.2	
Lym	18.5	24.0	
Mon	2.5	5.0	
Eos	1.3	1.9	

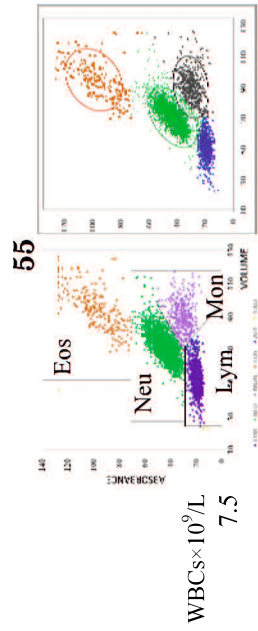


WBCs		SPC/ITC	
Manual	Pentra	Manual	SPC/ITC
Neu	72.0	75.5	
Lym	12.3	9.6	
Mon	12.7	11.0	
Eos	2.9	3.9	

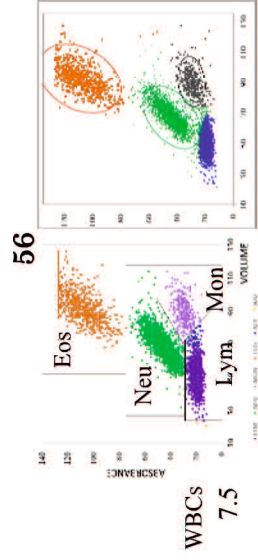


WBCs		SPC/ITC	
Manual	Pentra	Manual	SPC/ITC
Neu	65.8	71.5	
Lym	13.7	12.9	
Mon	6.2	5.2	
Eos	14.0	10.4	

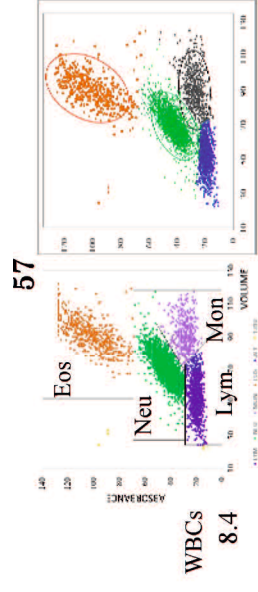
異常検体109名分の解析結果



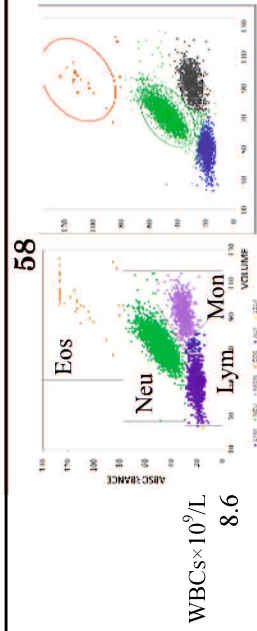
	Manual	SPC/ITC
WBCs×10 ⁹ /L	7.5	
Pentra	68.4	71.5
Neu	74.0	15.0
Lym	15.0	9.1
Mon	4.5	4.4
Eos	4.8	



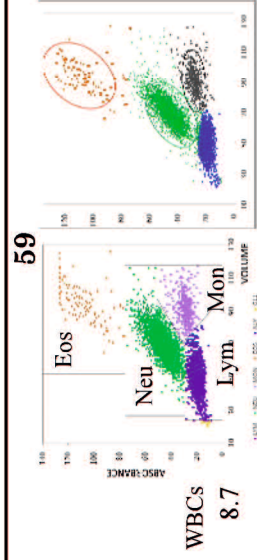
	Manual	SPC/ITC
WBCs	7.5	
Pentra	42.0	50.4
Neu	34.0	29.7
Lym	6.7	7.2
Mon	18.5	12.8



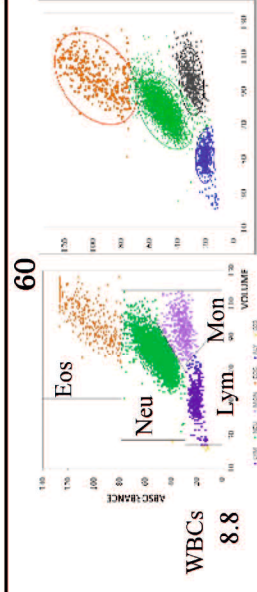
	Manual	SPC/ITC
WBCs	8.4	
Pentra	57.7	61.9
Neu	21.8	20.4
Lym	8.3	8.5
Mon	12.3	9.3



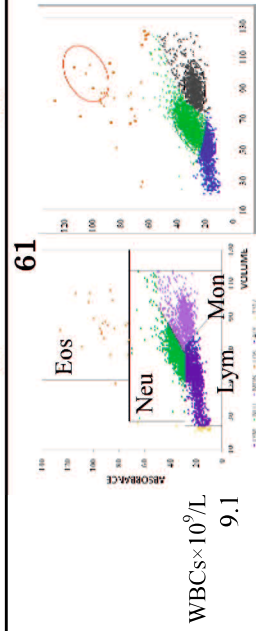
	Manual	SPC/ITC
WBCs×10 ⁹ /L	8.6	
Pentra	57.5	60.1
Neu	25.5	24.0
Lym	10.0	15.6
Mon	0.5	0.3
Eos	0.8	



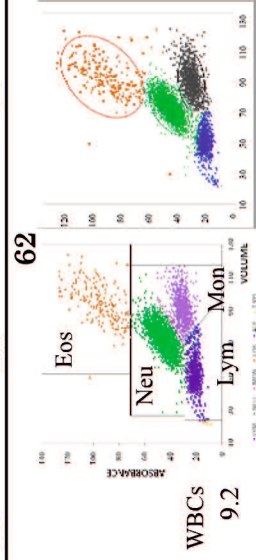
	Manual	SPC/ITC
WBCs	8.7	
Pentra	56.3	58.9
Neu	33.7	31.9
Lym	7.6	7.5
Mon	2.0	1.7



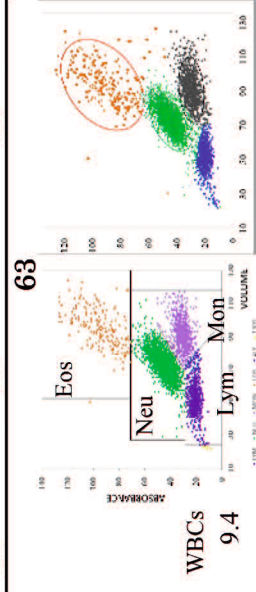
	Manual	SPC/ITC
WBCs	8.8	
Pentra	76.0	78.0
Neu	8.1	7.7
Lym	8.6	7.4
Mon	6.0	6.9



	Manual	SPC/ITC
WBCs×10 ⁹ /L	9.1	
Pentra	39.4	61.3
Neu	28.8	19.7
Lym	5.8	18.6
Mon	1.8	0.4

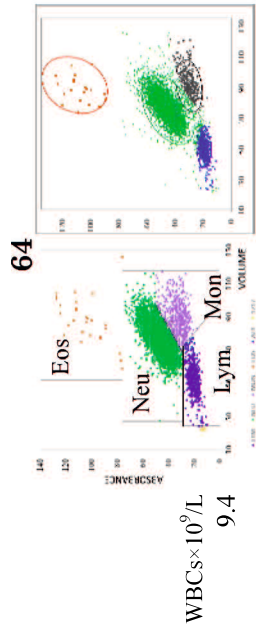


	Manual	SPC/ITC
WBCs	9.2	
Pentra	61.9	64.4
Neu	14.3	14.4
Lym	6.3	17.0
Mon	4.8	4.2

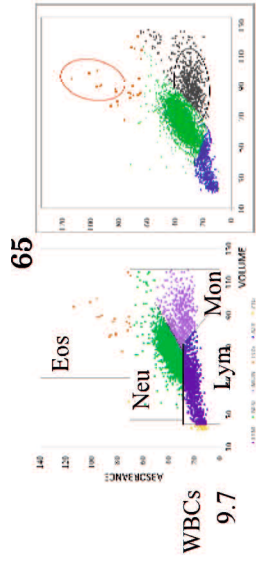


	Manual	SPC/ITC
WBCs	9.4	
Pentra	71.4	64.4
Neu	22.3	14.3
Lym	5.8	17.0
Mon	2.8	4.2

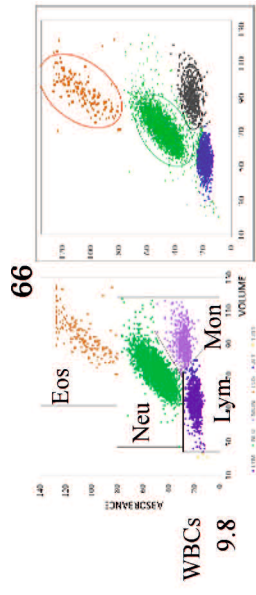
異常検体109名分の解析結果



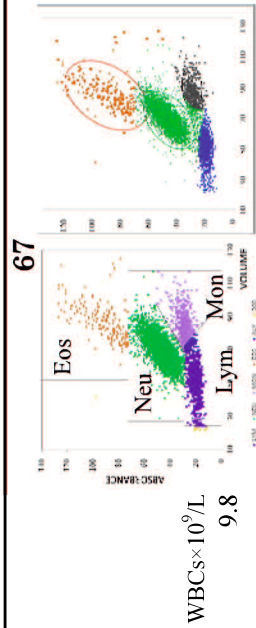
	Manual	SPC/ITC
Neu	80.8	86.2
Lym	14.5	9.7
Mon	3.3	3.8
Eos	1.0	0.4



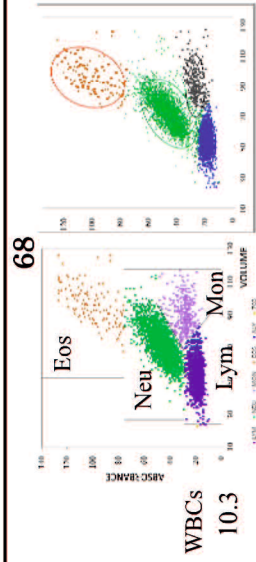
	Manual	SPC/ITC
Neu	67.5	79.8
Lym	20.8	10.9
Mon	10.9	8.8
Eos	0.5	0.4



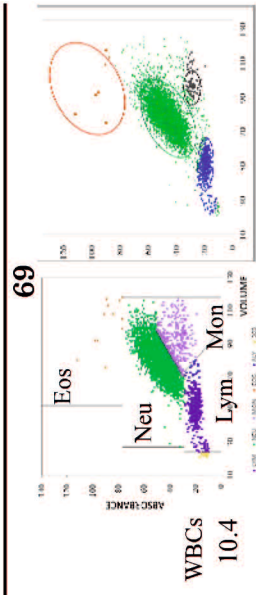
	Manual	SPC/ITC
Neu	69.1	73.1
Lym	18.5	16.3
Mon	9.1	8.3
Eos	3.0	2.3



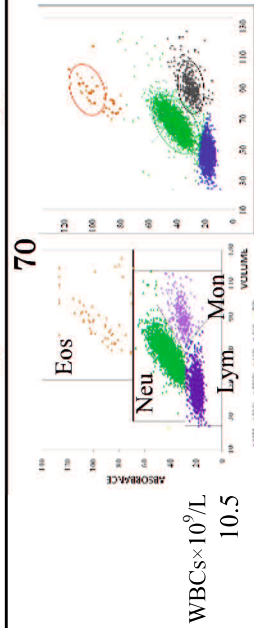
	Manual	SPC/ITC
Neu	74.1	78.0
Lym	14.7	11.9
Mon	8.4	7.0
Eos	2.7	3.1



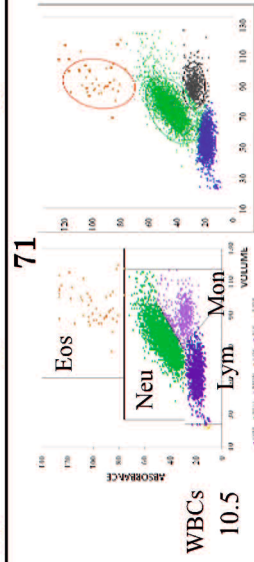
	Manual	SPC/ITC
Neu	60.9	64.9
Lym	31.0	28.8
Mon	5.7	4.4
Eos	1.9	1.9



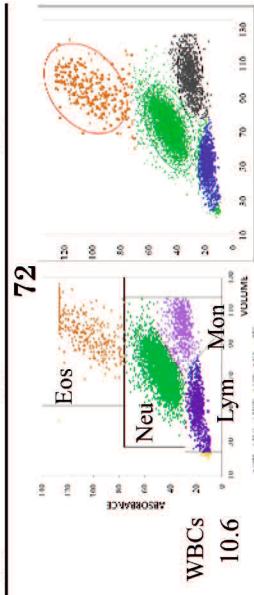
	Manual	SPC/ITC
Neu	88.4	93.5
Lym	6.7	5.7
Mon	4.5	0.7
Eos	0.3	0.1



	Manual	SPC/ITC
Neu	64.0	67.5
Lym	29.4	26.8
Mon	5.3	5.0
Eos	0.9	0.6

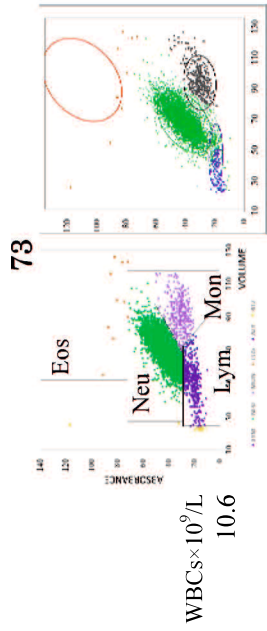


	Manual	SPC/ITC
Neu	62.1	68.6
Lym	26.0	23.6
Mon	10.3	7.0
Eos	0.9	0.8



	Manual	SPC/ITC
Neu	63.8	67.1
Lym	18.5	17.3
Mon	11.8	11.0
Eos	5.5	4.6

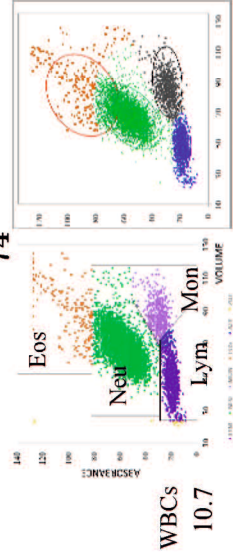
異常検体109名分の解析結果



73

WBCs×10⁹/L
10.6

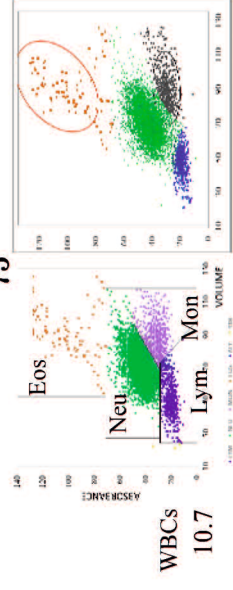
	Manual	SPC/ITC
Neu	90.8	93.8
Lym	4.3	2.1
Mon	4.8	4.1
Eos	0.3	0.0



74

WBCs
10.7

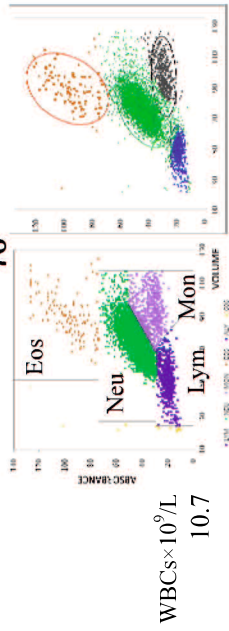
	Manual	SPC/ITC
Neu	73.5	70.6
Lym	18.3	17.9
Mon	5.0	8.7
Eos	2.5	2.5



75

WBCs
10.7

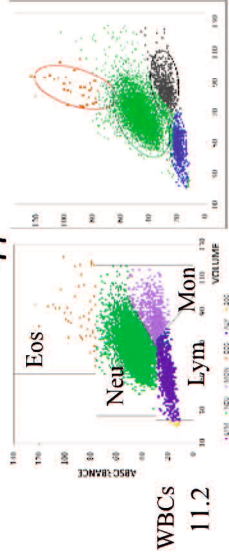
	Manual	SPC/ITC
Neu	87.3	79.8
Lym	7.0	8.3
Mon	4.0	10.0
Eos	1.8	1.5



76

WBCs×10⁹/L
10.7

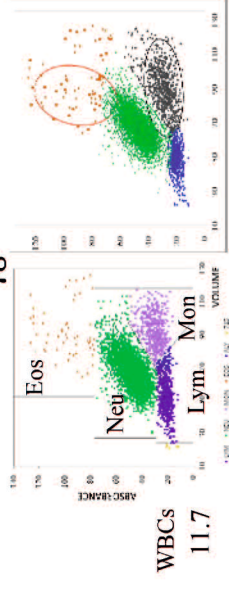
	Manual	SPC/ITC
Neu	84.0	85.9
Lym	8.3	7.1
Mon	3.0	5.3
Eos	1.8	1.7



77

WBCs
11.2

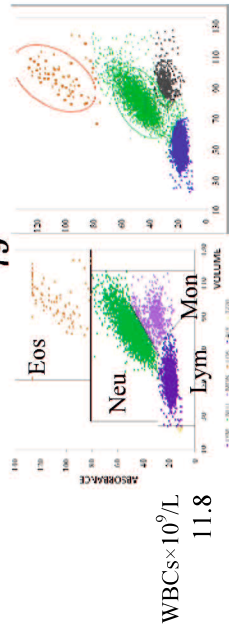
	Manual	SPC/ITC
Neu	77.8	76.6
Lym	10.8	11.6
Mon	10.0	10.4
Eos	1.3	0.8



78

WBCs
11.7

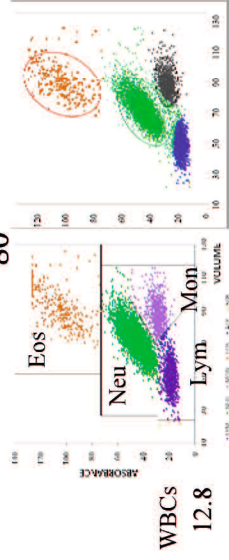
	Manual	SPC/ITC
Neu	88.5	84.6
Lym	6.8	7.9
Mon	4.0	6.7
Eos	0.5	0.7



79

WBCs×10⁹/L
11.8

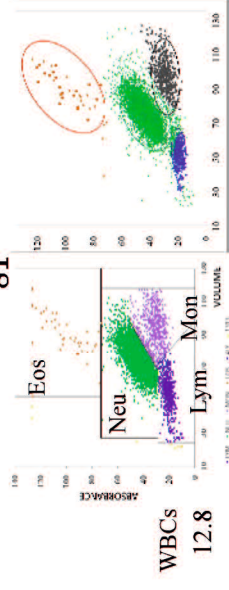
	Manual	SPC/ITC
Neu	33.0	56.2
Lym	58.8	37.9
Mon	4.8	5.0
Eos	1.3	0.9



80

WBCs
12.8

	Manual	SPC/ITC
Neu	65.8	66.0
Lym	23.8	19.9
Mon	6.8	9.7
Eos	3.0	4.0

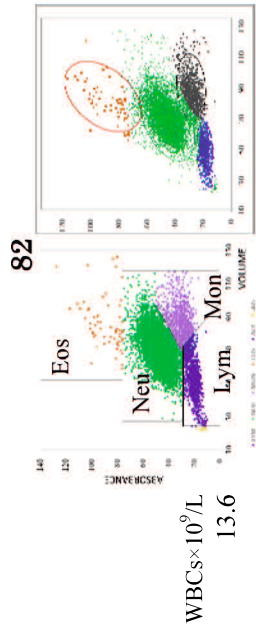


81

WBCs
12.8

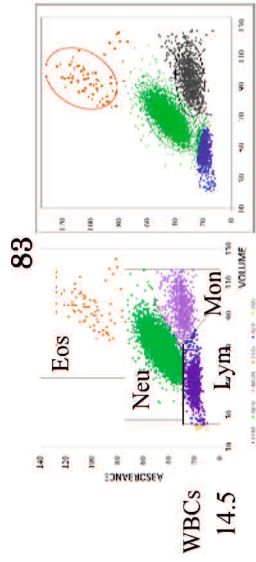
	Manual	SPC/ITC
Neu	89.0	81.3
Lym	7.0	9.2
Mon	3.3	8.7
Eos	0.3	0.6

異常検体109名分の解析結果



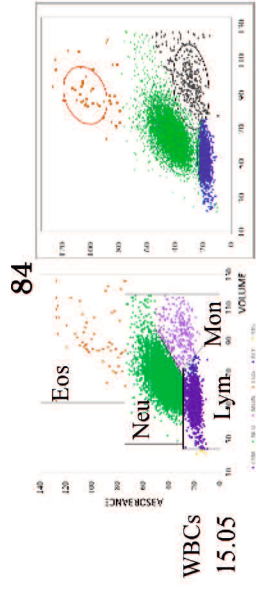
WBCs×10⁹/L
13.6

	Manual	SPC/ITC
Neu	78.8	84.2
Lym	9.5	7.7
Mon	10.8	7.6
Eos	1.0	0.5



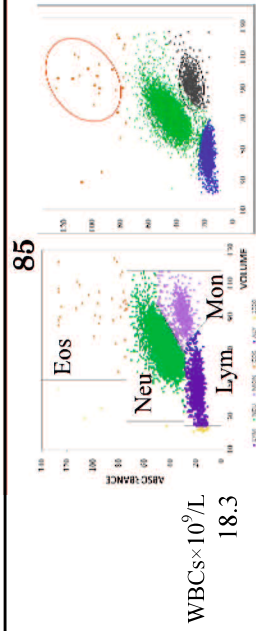
WBCs
14.5

	Manual	SPC/ITC
Neu	80.0	76.3
Lym	12.3	12.3
Mon	6.3	10.0
Eos	1.0	0.9



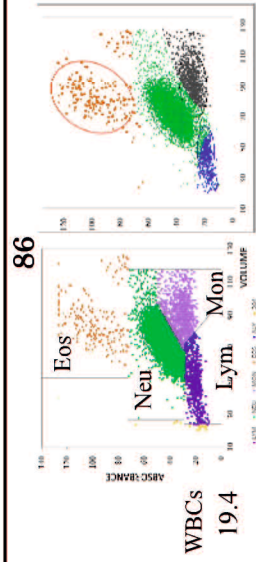
WBCs
15.05

	Manual	SPC/ITC
Neu	82.0	79.2
Lym	16.0	15.9
Mon	0.8	3.7
Eos	0.8	0.6



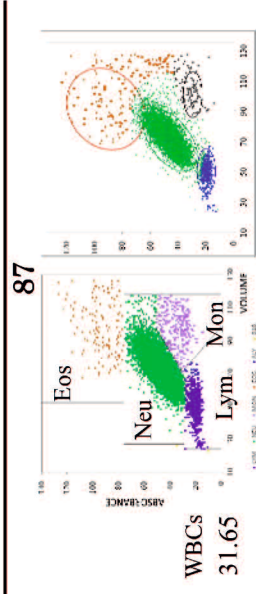
WBCs×10⁹/L
18.3

	Manual	SPC/ITC
Neu	71.3	76.3
Lym	24.8	22.6
Mon	3.3	1.0
Eos	0.3	0.0



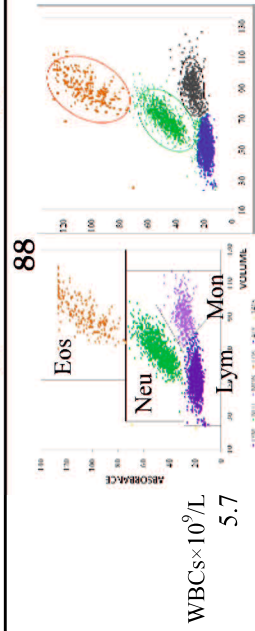
WBCs
19.4

	Manual	SPC/ITC
Neu	84.3	79.9
Lym	7.3	7.0
Mon	5.3	10.4
Eos	3.0	2.0



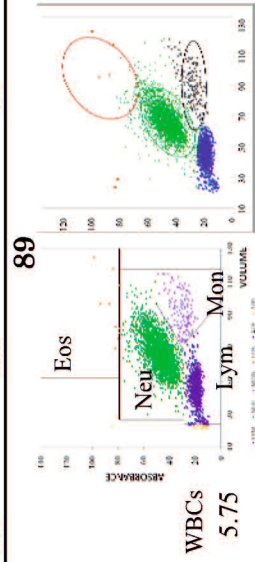
WBCs
31.65

	Manual	SPC/ITC
Neu	95.8	92.9
Lym	3.5	4.1
Mon	0.5	1.9
Eos	0.3	0.8



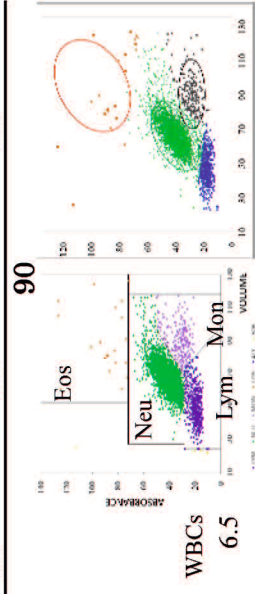
WBCs×10⁹/L
5.7

	Manual	SPC/ITC
Neu	35.8	36.2
Lym	49.0	46.0
Mon	5.5	12.5
Eos	8.0	5.3



WBCs
5.75

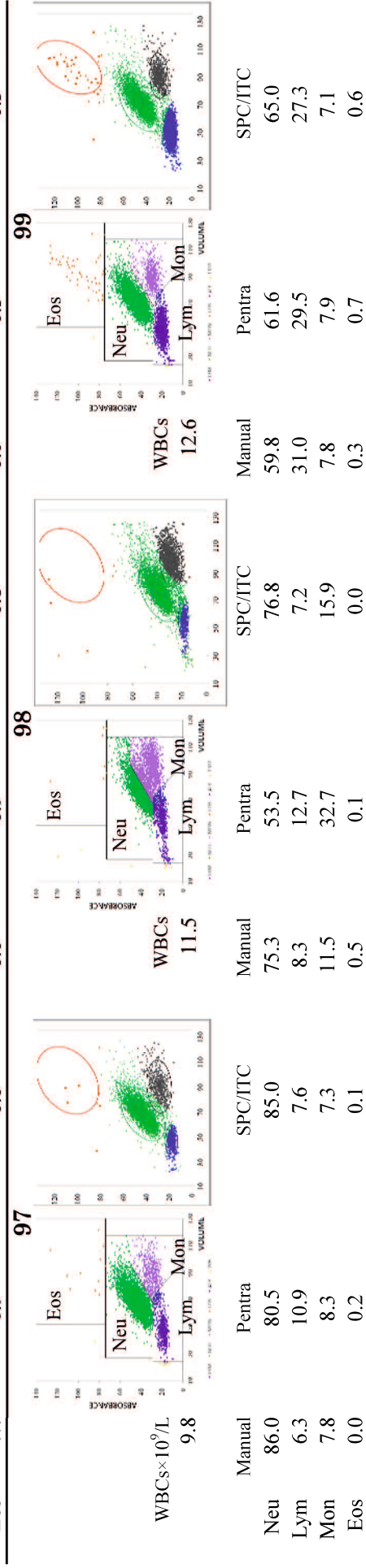
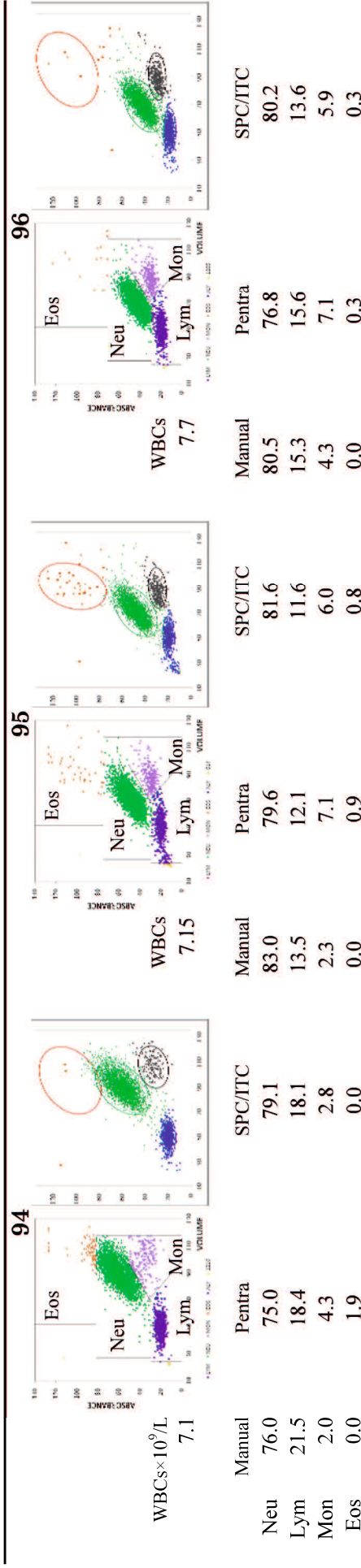
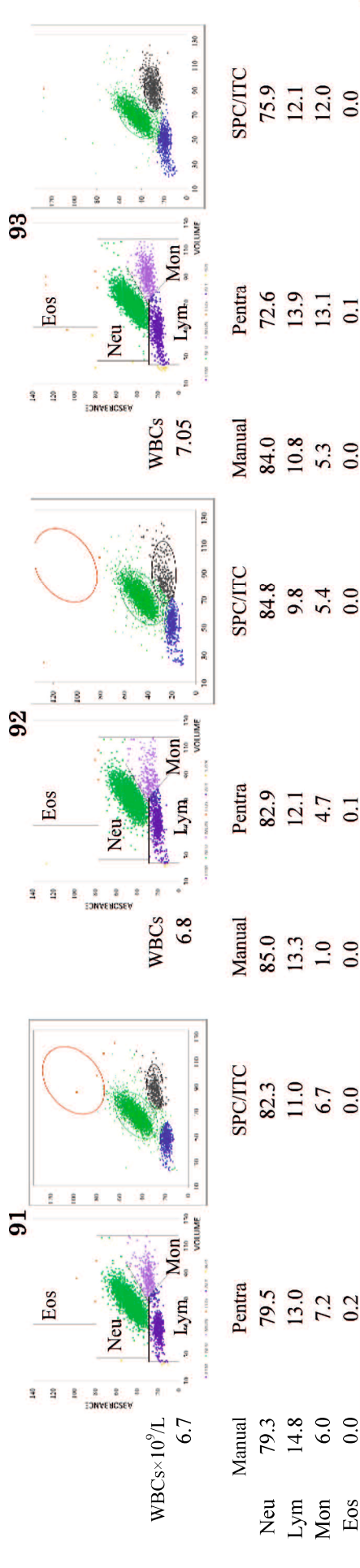
	Manual	SPC/ITC
Neu	75.8	73.1
Lym	22.0	22.7
Mon	2.3	3.4
Eos	0.0	0.3



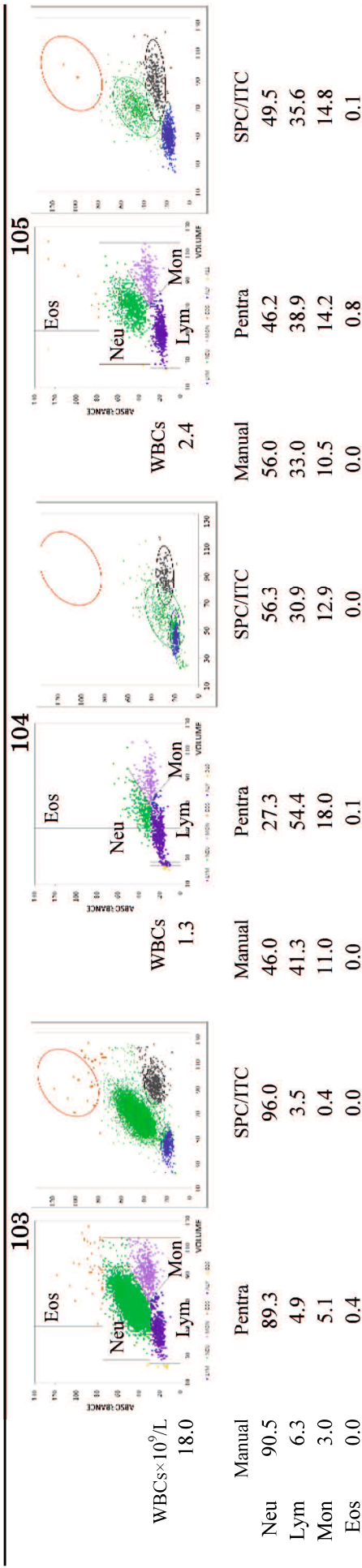
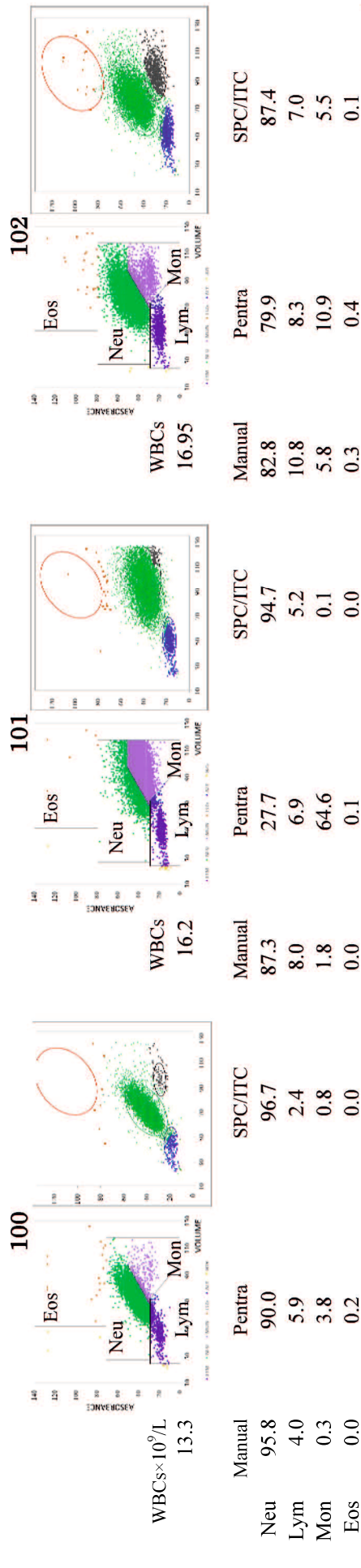
WBCs
6.5

	Manual	SPC/ITC
Neu	87.5	79.2
Lym	7.5	13.4
Mon	3.3	6.5
Eos	0.3	0.5

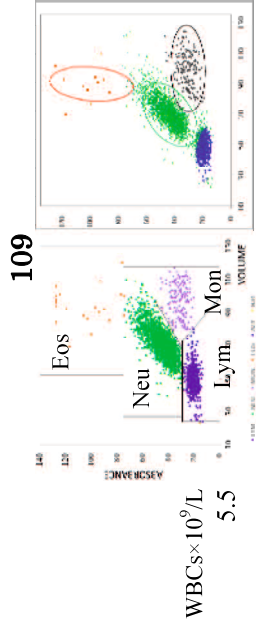
異常検体109名分の解析結果



異常検体109名分の解析結果



異常検体109名分の解析結果



	Manual	Pentra	SPC/ITC
Neu	80.0	72.9	77.2
Lym	17.0	21.5	17.9
Mon	2.8	4.8	4.5
Eos	0.0	0.7	0.4