

(19) 日本国特許庁(JP)

(12) 公表特許公報(A)

(11) 特許出願公表番号

特表2016-511397

(P2016-511397A)

(43) 公表日 平成28年4月14日(2016.4.14)

(51) Int.Cl.	F I	テーマコード (参考)
GO 1 N 33/48 (2006.01)	GO 1 N 33/48 M	2 GO 4 5
GO 1 N 33/53 (2006.01)	GO 1 N 33/53 Y	
GO 1 N 33/543 (2006.01)	GO 1 N 33/53 K	
GO 1 N 33/536 (2006.01)	GO 1 N 33/543 5 9 7	
	GO 1 N 33/536 D	

審査請求 未請求 予備審査請求 未請求 (全 44 頁)

(21) 出願番号 特願2015-555731 (P2015-555731)
 (86) (22) 出願日 平成26年1月31日 (2014.1.31)
 (85) 翻訳文提出日 平成27年9月4日 (2015.9.4)
 (86) 国際出願番号 PCT/EP2014/051963
 (87) 国際公開番号 WO2014/118343
 (87) 国際公開日 平成26年8月7日 (2014.8.7)
 (31) 優先権主張番号 13153512.2
 (32) 優先日 平成25年1月31日 (2013.1.31)
 (33) 優先権主張国 欧州特許庁 (EP)

(71) 出願人 515121759
 ユニベルシテ ドゥ モンペリエ
 フランス国, 34090 モンペリエ, リ
 ュ オーギュスト ブルソネ 163
 (71) 出願人 515208625
 サントル オスピタリエ ユニベルシテイ
 ル ドゥ モンペリエ
 CENTRE HOSPITALIER
 UNIVERSITAIRE DE MO
 NTPPELLIER
 フランス共和国, エフ-34295 モン
 ペリエ セデックス 5, アヴニユ デュ
 ドワイヤン ガストン ジロー 191

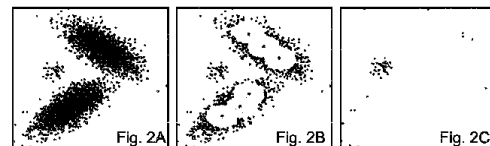
最終頁に続く

(54) 【発明の名称】 希少事象を同定する方法

(57) 【要約】

本発明は、細胞大集団中の特異的細胞の亜集団を同定するための方法において、前記大集団の細胞を n 種の試薬に曝露するステップと、前記 n 種の試薬を検出するステップと、クラスタ化により細胞を k 個の異なるクラスタにグループ化するステップと、希少細胞でない細胞を除去するステップと、を含む方法に関する。

【選択図】 図 2



Figures 2

【特許請求の範囲】

【請求項 1】

n次元空間内で細胞大集団中の特異的細胞の亜集団を同定するための方法において、

a. 前記大集団の細胞をn種の試薬に曝露するステップであって、前記n種の試薬が、前記大集団の各細胞のn種の異なる構成要素の存在、不存在または量の検出を可能にするステップであって、nは2以上であるステップと；

b. 各細胞をn次元空間内部の特定の位置に割当てするために、前記大集団に属する各細胞について前記n種の試薬を検出するステップと；

c. 細胞をクラスタ化によってk個の異なるクラスタにグループ化するステップであって、各クラスタが中心 C_k および半径Dにより特徴づけされており、クラスタ化は、前記大集団に属する細胞の20%~90%が前記k個のクラスタの1つに割当てられるようなものであり、kおよび C_k のパラメータは、前記規定のクラスタに割当てられる細胞の前記百分率に左右され、ここでクラスタ化ステップが、k平均修正アルゴリズムを実施することによって達成されるステップと；

d. 隣接するクラスタをグループ化してより大きいクラスタを得、ここで隣接するクラスタは、2つのクラスタの中心C間のユークリッド距離が半径Dの2倍よりも小さいようなものであり、かつより大きい前記クラスタの中心 C_{1k} ならびにより大きい前記クラスタに属する細胞の共分散行列を推定するステップと；

e. 前記n次元の各々の中より大きいクラスタの半径を0.01から0.1まで変動する係数だけ増加させることによって各々の拡大クラスタについてのすべり領域を定義づけし、前記すべり領域に属する各細胞について、マハラノビス距離を計算するステップと；

f. $D_{1k}(1 + \quad)$ よりも小さいマハラノビス距離を有する1セットの細胞に属する細胞の数を推定し、前記セットの密度を測定するステップであって、前記セットの細胞が、すべり領域に属するがより大きいクラスタには属していない細胞に対応しており、したがって、

- 密度が、10超である値Nより高く、好ましくはNが10から1000まで変動し、詳細にはNが10から500まで変動する場合、前記セットは前記特異的細胞を含まないものとみなされ、前記セットの密度が前記値Nよりも低くなるまでステップfおよびgがp回反復され、前記セットは $D_{1k}(1 + \quad)$ pより小さいマハラノビス距離を有する細胞により定義され、

- 前記セットの密度がN以下である場合、前記セットは前記特異的細胞を含んでいる、

ステップと；

を含む方法。

【請求項 2】

前記n種の試薬が、細胞タンパク質、脂質、糖質または核酸分子と相互作用する蛍光試薬である、請求項1に記載の方法。

【請求項 3】

検出ステップbがフローサイトメトリーによって実施される、請求項1または2に記載の方法。

【請求項 4】

前記細胞大集団が、血液標本、脳脊髄液、羊水、気管支肺胞洗浄液、母乳、頸腔部液などのあらゆる動物またはヒトの体液に由来するものである、請求項1~3のいずれか一項に記載の方法。

【請求項 5】

血液標本中の成熟内皮細胞の亜集団を同定することを目的とし、大集団の細胞が、少なくともCD45、CD105およびCD146というマーカーで標識されている、請求項1~4のいずれか一項に記載の方法。

【請求項 6】

10

20

30

40

50

血液標本中の前駆内皮細胞の亜集団を同定することを目的とし、大集団の細胞が、少なくともCD45、CD34、CD133、およびCD309というマーカーで標識されている、請求項1～4のいずれか一項に記載の方法。

【請求項7】

上皮細胞の亜集団を同定することを目的とし、大集団の細胞が少なくともCD326、CD45、およびサイトケラチンに向けられた抗体といったマーカーで標識されている、請求項1～4のいずれか一項に記載の方法。

【請求項8】

調節B細胞の亜集団またはエプスタイン・バールウイルス（EBV）に感染した記憶B細胞の亜集団を同定することを目的とし、大集団の細胞が少なくともCD27、CD24、CD19、およびIL-10というマーカー、またはCD27、CD19およびEBV抗原に向けられた抗体というマーカーでそれぞれ標識されている、請求項1～4のいずれか一項に記載の方法。

10

【請求項9】

調節T細胞の亜集団を同定することを目的とし、大集団の細胞が、少なくともCD4、CD25、Foxp3およびサイトカインに向けられた抗体というマーカーで標識されている、請求項1～4のいずれか一項に記載の方法。

【請求項10】

ヒト免疫不全ウイルス（HIV）に感染したCD4+T細胞の亜集団を同定することを目的とし、大集団の細胞が少なくともCD4、CD3、CD25およびHIV抗原に向けられた抗体というマーカーで標識されている、請求項1～4のいずれか一項に記載の方法。

20

【請求項11】

= 10 - 1であり、N = 10である、請求項1～10のいずれか一項に記載の方法。

【請求項12】

癌、血管および免疫病理および感染性疾患を含む病理のインビトロ診断またはインビトロ予後診断のための方法において、請求項1～11のいずれか一項に記載のn次元空間内の細胞大集団中の特異的細胞の亜集団を同定するステップを含む方法。

【請求項13】

請求項1～12のいずれか一項に記載の方法のステップc～fを実施できるようにする適切な媒体上のコンピュータプログラム。

30

【請求項14】

a. 前記大集団の各細胞のn種の異なる構成要素の存在、不存在または量を検出するためのn種の試薬であって、nは2以上である、試薬と、
b. 前記n種の異なる細胞マーカーの検出用手段と、
c. 請求項13に記載のコンピュータプログラムと、
を含むキット。

【発明の詳細な説明】

【技術分野】

【0001】

本発明は、希少事象を同定するための方法に関する。

40

【背景技術】

【0002】

不均一細胞集団を特徴づけするための1つの方法は、フローサイトメトリーによる方法である。この技術を用いると、細胞は、染料に接合した抗体を用いて標識される。フローサイトメトリーは、日常的に1個、2個またはそれ以上の免疫蛍光マーカーを同時に定量的に検出することができる。多数の免疫蛍光標識を細胞の光散乱特性と組み合わせることによって、異なる系統の細胞間のみならず、これらの系統内のさまざまな成熟段階にある細胞の間で区別を行なうことが可能となる。フローサイトメトリーにより同定された集団は、次に、計器上で利用可能な細胞選別用電子機器を用いて単離される。

50

【 0 0 0 3 】

国際公開第 2 0 0 6 / 0 8 9 1 9 0 号は、多次元解析を用いて異常細胞を検出するための方法を開示している。この出願は、規定の標本の細胞をクラスタ化し基準標本と比較して、異常細胞、すなわち前記規定の標本内にのみ存在する希少細胞を同定する方法を開示している。

【 0 0 0 4 】

この方法は、大集団に属する希少細胞の同定を可能にするものの、正常な細胞集団の使用を必要とする。

【 0 0 0 5 】

個体または個体群について前記対照標本の獲得が不可能である場合、求められる比較ステップの実施は困難であり得る。

【 発明の概要 】

【 発明が解決しようとする課題 】

【 0 0 0 6 】

したがって、いかなる標本であれ、細胞大集団内の細胞の亜集団の同定を可能にし、かつ基準標本に依存しない方法を提供する必要性が存在する。

【 0 0 0 7 】

本発明の目的は、上述の不都合を克服することにある。

【 0 0 0 8 】

本発明の別の目的は、技術者の介入とは無関係に直接的で単純かつ再現性ある形で希少細胞の検出を可能にする方法を提供することにある。

【 0 0 0 9 】

本発明のさらに別の目的は、上述の方法を実施することのできるコンピュータプログラムを提供することにある。

【 課題を解決するための手段 】

【 0 0 1 0 】

本発明は、 n 次元空間内で細胞大集団中の特異的細胞の亜集団を同定するための方法において、

a. 前記大集団の細胞を n 種の試薬に曝露するステップであって、前記 n 種の試薬が、前記大集団の各細胞の n 種の異なる構成要素の存在、不存在または量の検出を可能にするステップであって、 n は 2 以上であるステップと；

b. 各細胞を n 次元空間内部の特定の位置に割当てするために、前記大集団に属する各細胞について前記 n 種の試薬を検出するステップと；

c. 細胞をクラスタ化によって k 個の異なるクラスタにグループ化するステップであって、各クラスタが中心 C_k および半径 D により特徴づけされており、クラスタ化は、前記大集団に属する細胞の 20% ~ 90% が前記 k 個のクラスタの 1 つに割当てられるようなものであり、 k および C_k のパラメータは、前記規定のクラスタに割当てられる細胞の前記百分率に左右されるステップと；

d. 隣接するクラスタをグループ化してより大きいクラスタを得、ここで隣接するクラスタは、2つのクラスタの中心 C_k 間のユークリッド距離が半径 D の 2 倍よりも小さいようなものであり、かつより大きい前記クラスタの中心 C_{1k} ならびにより大きい前記クラスタに属する細胞の共分散行列を推定するステップ、ここでより大きい前記クラスタは半径 D_{1k} を有するステップと；

e. 前記 n 次元の各々の中より大きいクラスタの半径を 0.01 から 0.1 まで変動する係数だけ増加させることによって各々の拡大クラスタについてのすべり領域を定義付けし、前記すべり領域に属する各細胞について、マハラノビス距離を計算するステップと；

f. $D_{1k} (1 + \quad)$ よりも小さいマハラノビス距離を有する 1 セットの細胞に属する細胞の数を推定し、前記セットの密度を測定するステップであって、前記セットの細胞が、すべり領域に属するがより大きいクラスタには属していない細胞に対応しており、したがっ

10

20

30

40

50

て、

- 密度が、10超である値Nより高く、好ましくはNが10から1000まで変動し、詳細にはNが10から500まで変動する場合、前記セットは前記特異的細胞を含まないものとみなされ、前記セットの密度が前記値Nよりも低くなるまでステップfおよびgがp回反復され、前記セットは $D_{1k} (1 +)^p$ より小さいマハラノビス距離を有する細胞により定義され、

- 前記セットの密度がN以下である場合、前記セットは前記特異的細胞を含んでいる、
ステップと；

を含む方法に関する。

【0011】

有利には、本発明は、n次元空間内で細胞大集団中の特異的細胞の亜集団を同定するための方法において、

a. 前記大集団の細胞をn種の試薬に曝露するステップであって、前記n種の試薬が、前記大集団の各細胞のn種の異なる構成要素の存在、不存在または量の検出を可能にするステップであって、nは2以上であるステップと；

b. 各細胞をn次元空間内部の特定の位置に割当てのために、前記大集団に属する各細胞について前記n種の試薬を検出するステップと；

c. 細胞をクラスタ化によってk個の異なるクラスタにグループ化するステップであって、各クラスタが中心 C_k および半径Dにより特徴づけされており、クラスタ化は、前記大集団に属する細胞の20%~90%が前記k個のクラスタの1つに割当てられるようなものであり、kおよび C_k のパラメータは、前記規定のクラスタに割当てられる細胞の前記百分率に左右され、ここでクラスタ化ステップが、k平均修正アルゴリズムを実施することによって達成されるステップと；

d. 隣接するクラスタをグループ化してより大きいクラスタを得、ここで隣接するクラスタは、2つのクラスタの中心C間のユークリッド距離が半径Dの2倍よりも小さいようなものであり、かつより大きい前記クラスタの中心 C_{1k} ならびにより大きい前記クラスタに属する細胞の共分散行列を推定するステップと；

e. 前記n次元の各々の中より大きいクラスタの半径を0.01から0.1まで変動する係数だけ増加させることによって各々の拡大クラスタについてのすべり領域を定義づけし、前記すべり領域に属する各細胞について、マハラノビス距離を計算するステップと；

f. $D_{1k} (1 +)$ よりも小さいマハラノビス距離を有する1セットの細胞に属する細胞の数を推定し、前記セットの密度を測定するステップであって、前記セットの細胞が、すべり領域に属するがより大きいクラスタには属していない細胞に対応しており、したがって、

- 密度が、10超である値Nより高く、好ましくはNが10から1000まで変動し、詳細にはNが10から500まで変動する場合、前記セットは前記特異的細胞を含まないものとみなされ、前記セットの密度が前記値Nよりも低くなるまでステップfおよびgがp回反復され、前記セットは $D_{1k} (1 +)^p$ より小さいマハラノビス距離を有する細胞により定義され、

- 前記セットの密度がN以下である場合、前記セットは前記特異的細胞を含んでいる、
ステップと；

を含む方法に関する。

【0012】

本発明は、細胞大集団内で多重データ解析を適用することで、前記大集団内に属する細胞小集団の同定が可能となるという本発明者らによる予想外の観察事実に基づいている。

【0013】

上述の通り、本発明に係る方法は、以下のステップを含む：

- 細胞を標識し、標識を同定する第1のステップ、

10

20

30

40

50

- 細胞をクラスタ化し；かつクラスタ化ステップを調整する第2のステップ、および
- クラスタ化を調整して、集団の細胞の大部分を含む最も正確なクラスタを得るステップであって、前記クラスタ化された細胞が除去されて希少細胞を明らかにするステップ。

【0014】

第1のステップ：細胞の標識および標識された細胞の検出

本発明に係る方法を実施するため、解析対象の集団に属する全ての細胞は、 n 種の試薬で標識され、ここで n は2以上である。

【0015】

細胞大集団に属する全ての細胞は、 n 種の試薬好ましくは、各々前記集団に属する各細胞の n 種の異なる構成要素と反応する n 種の異なる試薬で標識される。

【0016】

本発明においては、 $n \geq 2$ 、すなわち n は2、3、4、5、6、7、8、9、10、11、12、13、14、15、16、17、18、19、20またはそれ以上である。試薬の数は、細胞を判別するため1超でなければならず、試薬の数は、それらを同時にまたは逐次的に検出する施術者の能力によってのみ左右されるものと考えられる。

【0017】

n 種の試薬は、各細胞の構成要素に特異的であり、前記構成要素と試薬の反応が各細胞内の前記構成要素の存在、不存在または量を規定する。

【0018】

大集団の細胞を標識するために使用される n 種の試薬は検出可能である。検出は、それらの性質または物理的特性または化学的特性またはその両方に応じて特異的手段によって実施可能である。例えば、本発明の範囲を限定することなく、試薬は蛍光性、磁気性、リン光性、放射性、非水溶性、可活性化、可誘発性試薬であり得る。

【0019】

一般に、 n 種の試薬は、細胞の1つの特定の構成要素を特異的に認識する抗体である。検出可能となるために、抗体は、蛍光性染料、ビーズ、特に磁気ビーズ、酵素などの検出可能な化合物と接続される。

【0020】

試薬は同様に、DNAまたは任意の他の分子の挿入剤でもあり得る。

【0021】

集団の細胞と反応した試薬が検出され、こうして各細胞は n 種の特異的検出により同定され、これにより、 n 次元空間内で特定の位置に各細胞を割当てることができるようになる。存在、不存在または量は、規定の n 次元内の座標を決定する。

【0022】

第2のステップ：クラスタ化およびグループ化

各細胞が前記 n 次元空間内の特定の位置に割当てられた時点で、クラスタ化ステップが実施される。

【0023】

クラスタ解析またはクラスタリングは、同じグループ（クラスタと呼ばれる）内のオブジェクトが他のグループ（クラスタ）内のオブジェクトに比べて互いに（何らかの意味で）より類似しているような形で1セットのオブジェクトをグループ化するタスクである。それは探索的データマイニングの主要なタスクであり、多くの分野で使用される統計的データ解析用の一般的技術である。

【0024】

クラスタ解析自体は1つの特定のアルゴリズムではなく、解決すべき一般的タスクである。それは、1つのクラスタを構成しているものが何であるかそしていかにしてそれらを効率良く発見するかの概念において著しく異なっているさまざまなアルゴリズムにより達成可能である。一般的なクラスタの概念には、次の特質を伴うグループが含まれる：

- クラスタ構成員間の小さい距離、
- データ空間の高密度部域、

10

20

30

40

50

- 間隔または特定の統計的分布。

【0025】

したがって、クラスタリングは、多目的最適化問題として策定可能である。適切なクラスタリングアルゴリズムおよびパラメータ設定値（使用すべき距離関数、密度閾値または期待されるクラスタの数などの値を含む）は、個別のデータセットおよび結果の意図された用途によって左右される。クラスタ解析はそれ自体自動タスクではなく、試行錯誤が関与する知識発見または対話型多目的最適化の反復的プロセスである。多くの場合、結果が所望の特性を達成するまで前処理およびパラメータを修正することが必要となる。

【0026】

「クラスタ」の概念を精確に定義づけることは不可能であり、極めて多くのクラスタリングアルゴリズムが存在する理由の1つはここにある。当然のことながら、共通分母、すなわちデータオブジェクトグループは存在する。しかしながら、異なるクラスタモデルが使用され、これらのクラスタモデルの各々について再び異なるアルゴリズムが提供され得る。異なるアルゴリズムによって発見されるクラスタの概念は、その特性において著しく変動し、これらの「クラスタモデル」を理解することが、さまざまなアルゴリズム間の差異を理解するための鍵となる。典型的なクラスタモデルには、以下のものが含まれる：

- 接続性モデル：例えば、階層的クラスタリングは、距離接続性に基づいてモデルを構築する。
- セントロイドモデル：例えばk平均アルゴリズムは、単一平均ベクトルにより各クラスタを表現する。
- 分布モデル：クラスタは、期待値最大化アルゴリズムによって使用される多変量正規分布などの統計的分布を用いてモデリングされる。
- 密度モデル：例えばDBSCANおよびOPTICSは、データ空間内の接続された高密度領域としてクラスタを定義する。
- 部分空間モデル：バイクラスタリング（コクラスタリングまたは2モードクラスタリングとしても公知）においては、クラスタは、クラスタ構成員および関連する属性の両方でモデリングされる。
- グループモデル：一部のアルゴリズムは（残念なことに）、その結果についての改良モデルを提供せず、グループ化情報を提供するだけである。
- グラフベースのモデル：クリーク、すなわちエッジによりサブセット内の2つ毎のノードが接続されるようなグラフ内のノードサブセットを、クラスタの原型的形態とみなすことができる。完全接続性要件の緩和（エッジの一部が欠落し得る）は、準クリークとして公知である。

【0027】

セントロイドベースのクラスタリングにおいては、クラスタは、必ずしもデータセットの一構成員でなくてよい中央ベクトルによって表現される。クラスタ数がkに定められた場合、k平均クラスタリングが、1つの最適化問題として、以下のような形式上の定義を提供する：すなわち クラスタ中心を発見し、オブジェクトを最も近いクラスタ中心に割当て、こうしてクラスタからの距離の2乗が最小化されるようにすること。

【0028】

最適化問題はそれ自体、NP困難なものとして公知であり、したがって、一般的なアプローチは、近似解のみを捜し求めることにある。極めて周知の近似手法は、多くの場合実際には「k平均アルゴリズム」と呼ばれるロイドアルゴリズムである。しかしながら、それは局所最適を発見するにすぎず、一般的には、異なるランダム初期化を用いて多数回実行される。k平均の変形形態には多くの場合、多数の実行のうちの最良のものを選択することのみならず、データセットの構成員にセントロイドを制限すること（k-メドイド）、中央値を選択すること（k-中央値クラスタリング）、初期中心をより低いランダム度で選択すること（k-平均++）、またはファジークラスタ割当てを可能にすること（ファジーc-平均）などの最適化が含まれる。

【0029】

10

20

30

40

50

大部分のk平均タイプのアルゴリズムでは、クラスタ数 - k - が予め規定されている必要があり、これがこれらのアルゴリズムの最大の欠点の1つとみなされている。その上、アルゴリズムは、おおよそ類似のサイズのクラスタの方を好むが、これは、それらが常に1つのオブジェクトを最も近いセントロイドに割当てることになるからである。こうして多くの場合、クラスタ間において境界線が不正確にカットされることになる(アルゴリズムはクラスタの境界線ではなく、クラスタ中心を最適化したのであるから、これは意外なことではない)。

【0030】

K平均には数多くの興味深い理論特性がある。一方では、それはデータ空間をボロノイ図として公知の構造へと分割する。他方では、それは概念的に最近傍分類に近く、そのため機械学習において一般的である。第3に、それをモデルベースの分類の一変形形態として考えることができ、ロイドアルゴリズムを以下で論述するこのモデルのための期待値最大化アルゴリズムの一変形形態と考えることができる。

【0031】

有利には、本発明の中でk平均を使用する場合、kクラスタの各々は、以下のものにより定義される：

- セントロイド C_k および
- 半径 D 。

【0032】

換言すると、 $k = 3$ の場合、クラスタ1は C_1 および D により定義され、クラスタ2は C_2 および D により定義され、クラスタ3は、 C_3 および D により定義される。全てのクラスタは同じ半径を有するもののセントロイドは異なる。

【0033】

平面図形または2次元形状 X の幾何学的中心または重心とも呼ばれるセントロイドは、線を中心として等しいモーメントの2つの部分に X を分割する全ての直線の交差点である。非公式には、それは X の全ての点の「平均」(算術平均)である。

【0034】

定義は、 n 次元空間内の任意のオブジェクト X にまで及ぶ。すなわち、そのセントロイドは、 X を等しいモーメントの2つの部分に分割する全ての超平面の交差点である。

【0035】

発明者らは有利にも、空間の高密度領域内に存在する点のみを発見しクラスタ化するように構想されたk平均の修正済変形形態である *Dense K Means* を使用した。当業者であれば、上記の定義を考慮して、本発明に係るクラスタを得ることができる。

【0036】

クラスタが決定された場合、隣接するクラスタはグループ化されて、より大きいクラスタが得られる。

【0037】

本発明において、隣接するクラスタは、2つのクラスタの中心 C_k 間のユークリッド距離が半径 D_k の2倍よりも小さいようなものである。数学では、点 p と q の間のユークリッド距離は、それらを接続する線分の長さである。一般に、 n 次元空間について、距離は

【数1】

$$d(p, q) = \sqrt{(p_1 - q_1)^2 + (p_2 - q_2)^2 + \dots + (p_n - q_n)^2}$$

である。

【0038】

隣接するクラスタが同定され、より大きいクラスタが決定された場合、より大きい前記クラスタの中心 C_{1k} が推定される。その上共分散行列も同様に推定される。

【0039】

(分散行列または分散共分散行列としても公知である) 共分散行列は、 i 、 j 位置にあるその要素が、ランダムベクトル(すなわちランダム変数のベクトル)の i 番目および j 番目の要素の間の共分散である行列である。ベクトルの各要素は、有限数の視察経験値または全てのランダム変数の理論的同時確率分布によって規定される有限数または無限数の潜在値のいずれかを伴うスカラーランダム変数である。

【0040】

本発明において、より大きいクラスタの各々はこうしてその新たに定義された中心および共分散行列によって定義される。

【0041】

第3のステップ: クラスタの画定および高密度クラスタの除去

より大きいクラスタの決定に加えて、 n 次元の各々において半径サイズを係数 α だけ増大させることにより、各クラスタについてすべり領域が定義される。係数 α は 0.01 から 1 まで、有利には 0.05 から 0.5 まで変動し、詳細には α は 0.1 である。

【0042】

次に、前記すべり領域に属する各細胞についてのマハラノビス距離が計算される。

【0043】

マハラノビス距離は、相関独立変数によって定義される多次元空間内のセントロイドからのケースの距離である(独立変数が無相関である場合、それは単なるユークリッド距離と同じである)。こうして、この尺度は、1つの観察結果が独立変数値との関係における外れ値であるか否かの標示を提供する。

【0044】

ランダム変数の2つの標本(x 、 y)間のマハラノビス距離は、

【数2】

$$d(x,y) = \sqrt{(x-y)^T \Sigma^{-1} (x-y)}$$

として定義され、ここで Σ^{-1} は共分散行列の逆数である。

【0045】

すべり領域内では、全ての細胞は、 $D_{1k} (1 + \alpha)$ よりも小さいマハラノビス距離を有しており、ここで D_{1k} はより大きいクラスタ k の半径である。

【0046】

すべり領域に属する問題の細胞は、より大きいクラスタに属していない。

【0047】

このステップは、より大きいクラスタのサイズを係数 α だけ漸進的に増大させて、より大きい前記クラスタに属するがより大きい前記クラスタの境界線にある全ての細胞を得るために意図されている。

【0048】

すべり領域が決定された時点で、前記すべり領域の密度が評価される。密度が値 N よりも高く(N は 10 超である)、好ましくは、 N が 10 から 1000 まで変動し、詳細には、 N が 0 から 500 まで変動する(すなわち、これはすべり領域の密度が既定値 N よりも大きいことを意味する)場合、すべり領域はより大きいクラスタに属するものとみなされる。

【0049】

有利には、密度 N は 10 から 1000 まで、詳細には 10 から 500 まで変動し、詳細には 10 に等しい。

【0050】

こうして、本発明に係る方法のステップ f および g は、 $D_{1k} (1 + \alpha)^p$ よりも小さいマハラノビス距離を有する細胞により決定されるすべり領域 p の密度が N よりも小さくなるまで、 p 回繰返される。この場合、これはすなわち、すべり領域 n の細胞が、(係数により p 回拡大された)より大きいクラスタに属していないことを意味する。

10

20

30

40

50

【0051】

この段階で、より大きいクラスタに属する全ての細胞は除去され、残った細胞の中には、特異的細胞が期待される亜集団、すなわち希少細胞が存在する。

【0052】

有利には、本発明は、前記n種の試薬が、細胞タンパク質、脂質、糖質または核酸分子と相互作用する蛍光試薬である、上記で定義した方法に関する。

【0053】

蛍光性化合物は、その化合物の特徴である一波長範囲全体にわたり光エネルギーを吸収する。この光の吸収により、蛍光性化合物内の電子はより高いエネルギーレベルに高められる。励起電子は急速にその基底状態まで崩壊し、光量子として余剰のエネルギーを放出する。このエネルギー遷移は、蛍光と呼ばれる。

10

【0054】

蛍光性化合物を励起させることのできる範囲は、その吸収スペクトルと称される。吸収遷移では、蛍光遷移で放出されるよりも多くのエネルギーが消費されることから、発光波長は吸収される波長よりも長くなる。特定の化合物についての発光波長範囲は、発光スペクトルと称される。

【0055】

上述の通り、前記n種の試薬は、有利には、タンパク質、脂質、糖質または核酸分子などの細胞の構成要素と直接的または間接的に相互作用する蛍光性化合物または分子である。これには同様に、糖タンパク質、リン脂質および修飾糖質も含まれる。

20

【0056】

DNAと直接相互作用する蛍光性分子としてDNAの挿入分子を使用することができる。DNAを標識するための蛍光マーカーストしては、一般にヨウ化プロピジウム(PI)または7-アミノアクチノマイシンD(7-AAD)が使用される。当業者であれば、同じ特性を有する任意の他の作用物質を容易に選択することができる。

【0057】

詳細にはそのFc領域で蛍光性化合物と接続された抗体も同様に、タンパク質、脂質、糖質、タンパク質および脂質の糖型を検出するために使用することができる。蛍光性染料としては、一般に以下の化合物が使用される：FITC(フルオレセインイソチオシアネート)、Alexa Fluor(登録商標)488、R-PE(R-フィコエリスリン)、PE-Texas Red、PE-Alexa Fluor(登録商標)610、PE-Cy5、PerCP-Cy5.5、PerCP-eFluor(登録商標)710、PE-Cy7、Alexa Fluor(登録商標)532、APC(アロフィコシアニン)、eFluor(登録商標)660、Alexa Fluor(登録商標)647、Alexa Fluor(登録商標)700、APC-eFluor(登録商標)780、eFluor(登録商標)450、eFluor(登録商標)605NC、eFluor(登録商標)625NC、eFluor(登録商標)650NC、Pacific Blue(登録商標)、Pacific Orange(登録商標)、Brilliant Violet(登録商標)421、Brilliant Violet(登録商標)510、Brilliant Violet(登録商標)570、Brilliant Violet(登録商標)605、Brilliant Violet(登録商標)650、Brilliant Violet(登録商標)711およびBrilliant Violet(登録商標)785。このリストは限定的なものではなく、当業者であれば、最も適切な蛍光染料を容易に選択することができる。

30

40

【0058】

有利には、本発明に係る方法の中で使用されるn種の試薬は、分化クラスタ(CD)マーカーと呼ばれる膜貫通タンパク質の存在、不存在または量を特異的に検出する。一部の染料が細胞内にあり得るという点を指摘しておかなければならない(細胞内CD)。

【0059】

分化抗原クラスタは、白血球上で主に発現された膜タンパク質である。少数が、内皮細

50

胞、赤血球および幹細胞上でも発現される。分化抗原のクラスタは一般に、細胞マーカーとして使用され、細胞の表面上にどんな分子が存在するかに基づいて、細胞を定義づけるようにする。例えば、2つの一般に使用されるCD分子はCD4とCD8であり、これらは概して、それぞれT-リンパ球、ヘルパーおよび細胞毒性T細胞の2つの異なる亜型のためのマーカーとして使用される。CD4は、HIVにより特異的に認識され接続されて、ウイルス感染およびCD4+T細胞の破壊を導く。CD4+およびCD8+T細胞の相対存在度は、HIV感染の進行を監視するのに使用されるゴールドマーカーである。分化クラスタ(CD)抗原の発現の検出は、心臓血管疾患および腫瘍などの一部の疾病における診断方法として開発されるものと考えられている。

【0060】

ヒトにおいては、公知のCDマーカーの大部分が以下のものである：CD1a、CD1b、CD1c、CD1d、CD1e、CD2、CD3、CD3、CD3、CD4、CD5、CD5L、CD6、CD7、CD8a、CD8b、CD9、CD10/ネプリライシン、CD11a、CD11b/インテグリンアルファM、CD11c、Cdw12、CD13/ANPEP、CD14、CD15、CD15s、CD15u、CD15su、CD16a/FcガンマRIIA、CD16b/FcガンマRIIB、CD16-2/FCGR4、CD17、CD18/インテグリンベータ2/TNFRSF3、CD19、CD20/MS4A1、CD21、CD22、CD23/FCER2、CD24、CD25/IL-2R/IL-2RA、CD26/DPP4、CD27/TNFRSF7、CD27L/CD70/TNFSF7、CD28、CD29、CD30/TNFRSF8、CD30L/CD153/TNFSF8、CD31/PECAM1、CD32/FcガンマRII、CD32a/FcガンマRIIA、CD32b(変異2)、CD32b(変異3)、CD33/シグレック-3、CD34、CD34T、CD35、CD36/SCARB3、CD36L1/SCARB1、CD36L2/LIMP-2/SCARB2、CD37、CD38、CD39、CD40/TNFRSF5、CD40L/CD154/TNFSF5、CD41、CD42a、CD42b、CD42c/GP1BB、CD42d、CD43、CD44、CD44R、CD45、CD45RA、CD45RB、CD45RC、CD45RO、CD46、CD47、CD48/SLAMF2、CD49a、CD49b、CD49c、CD49d/インテグリンアルファ4、CD49e/インテグリンアルファ5、CD49f、CD50/ICAM-3、CD51、CD52、CD53、CD54/ICAM-1、CD55/DAF、CD56/NCAM1、CD57、CD58、CD59、CD60a、CD60b、CD60c、CD61/インテグリンベータ3、CD62E/E-セレクトリン、CD62L/L-セレクトリン、CD62P/P-セレクトリン、CD63、CD64/FcガンマRI、CD65、CD65s、CD66a/CEACAM1、CD66b/CD67/CEACAM8、CD66c/CEACAM6、CD66d/CEACAM3、CD66e/CEACAM5、CD66f、CD68、CD69、CD70/CD27L/TNFSF7、CD71/TFRC、CD72、CD73/NT5E、CD74、CD75/ST6GAL1、CD75s、CD77、CD79a、CD79b、CD80/B7-1、CD81、CD82/KAI-1、CD83、CD84/SLAMF5、CD85a、CD85b、CD85c、CD85d、CD85e、CD85f、CD85g、CD85h、CD85i、CD85j、CD85k、CD85I、CD85m、CD86/B7-2、CD87/PLAUR、CD88、CD89/FCAR、CD90/THY-1、CD91/LRP1、CD92、CD93/C1qR、CD94、CD95/APO-1/TNFRSF6、CD95L/CD178/TNFSF6、CD96、CD97、CD98/SLC3A2、CD99、CD99L2、CD99R、CD100/SEMA4D、CD101、CD102/ICAM-2、CD103、CD104、CD105/エンドグリン、CD106/VCAM-1、CD107a/LAMP-1、CD107b/LAMP2、Cdw108、CD109、CD110/TPOR/C-MPL、CD111/ネクチン-1/PVRL1、CD112/ネクチン-2、CD113/ネクチン-3、CD114/G-CSFR、CD115/CSF1R/M

10

20

30

40

50

CSFレセプター、CD116/GM-CSFR、CD117/c-キット、CD118
 /LIFR、CD119/IFNGR1、CD120a/TNFR1/TNFRSF1A
 、CD120b/TNFR2/TNFRSF1B、CD121a/IL-1R1、CD1
 21b/IL-1R2、CD122/IL-2RB、CD123/IL-3RA、CD1
 24/IL-4R、CD125/IL-5RA、CD126/IL-6R、CD127/
 IL-7RA、CD128/CD181/CXCR1、CD128b/CD182/CX
 CR2、Cdw129、CD130/gp130/IL6ST、CD131/IL-3R
 B/CSF2RB、CD132/IL-2RG、CD133、CD134/OX40/T
 NFRSF4、CD135/FLT3/FLK2、CD136/MST1R、CD137
 /4-1BB/TNFRSF9、CD137L/4-1BBL/TNFSF9、CD13
 8/シンデカン-1/SDC1、CD139、CD140a/PDGFRA、CD140
 b/PDGF RB、CD141、CD142/組織因子、CD143、CD144/VE
 -カドヘリン、Cdw145、CD146/MCAM、CD147/EMMPRIN、C
 D148、CD150/SLAM、CD151、CD152/CTLA-4、CD153
 /CD30L/TNFSF8、CD154/CD40L/TNFSF5、CD155/P
 VR、CD155b、CD156a/ADAM8、CD156b、CD156c、CD1
 57/BST1、CD158a、CD158b1、CD158b2/KIR2DL3、C
 D158c、CD158d、CD158e1、CD158e2、CD158f、CD15
 8g、CD158h、CD158i、CD158j、CD158k、CD158z、CD
 159a、CD159c、CD160、CD161、CD162/PSGL-1、CD1
 62R、CD163、CD164、CD165、CD166/ALCAM、CD167a
 /DDR1/MCK10、CD168、CD169、CD170、CD171/NCAM
 -L1、CD172a/SIRPアルファ、CD172b/SIRPベータ、CD172
 g/SIRPガンマ、CD173、CD174、CD175、CD175s、CD176
 、CD177、CD178/CD95L/TNFSF6、CD179a、CD179b、
 CD180/RP105、CD181/CD128/CXCR1、CD182/CD12
 8b/CXCR2、CD183、CD184、CD185、CD186、CD191、C
 D192、CD193、CD194、CD195、CD196、CD197、Cdw19
 8、Cdw199、CD200、CD200R、CD200R1、CD200R4、CD
 200RLa、CD201、CD202b/Tie2、CD203c、CD204/MS
 R1、CD205、CD206、CD207/ランゲリン、CD208/DC-LAMP
 、CD209/DC-SIGN、CD209b/SIGNR1、CD209g、CD21
 0a/IL-10RA、Cdw210b/IL-10RB、CD212/IL12RB1
 、CD213a1/IL-13RA1、CD213a2/IL-13RA2、CD217
 /IL-17R/IL-17RA、CD218a/IL-18R1/IL-18RA、C
 D218b/IL-18RAP/IL-1R7、CD220/インシュリンR、CD22
 1/IGF1R、CD222、CD223、CD224、CD225、CD226/DN
 AM-1、CD227/MUC-1/ムチン1、CD228、CD229/LY9、CD
 230、CD231、CD232、CD233、CD234、CD235a、CD235
 ab、CD235b、CD236、CD236R、CD238、CD239/BCAM、
 CD240CE、CD240D、CD240DCE、CD241CD242、CD243
 、CD244/2B4/SLAMF4、CD245、CD246、CD247、CD24
 8、CD249/ENPEP、CD252、CD253/TNFSF10/TRAIL、
 CD254/RANKL/OPGL/TNFSF11、CD255、CD256/TNF
 SF13、CD257/BlyS/TNFSF13B、CD258/LIGHT/TNF
 SF14、CD261/TRAILR1/TNFRSF10A、CD262/TRAIL
 R2/TNFRSF10B、CD263/TRAILR3/TNFRSF10C、CD2
 64/TRAILR4/TNFRSF10D、CD265、CD266/TWEAKR/
 TNFRSF12A、CD267/TACI/TNFRSF13B、CD268/BAF
 FR/TNFRSF13C、CD269/TNFRSF17/BCMA、CD270、C

10

20

30

40

50

D 2 7 1、C D 2 7 2、C D 2 7 3 / B 7 - D C / P D - L 2、C D 2 7 4 / B 7 - H 1 / P D - L 1、C D 2 7 5、C D 2 7 6 / B 7 - H 3、C D 2 7 7、C D 2 7 8 / I C O S / A I L I M、C D 2 7 9 / P D 1 / P D C D 1、C D 2 8 0、C D 2 8 1 / T L R 1、C D 2 8 2 / T L R 2、C D 2 8 3 / T L R 3、C D 2 8 4 / T L R 4、C D 2 8 6、C D 2 8 8、C D 2 8 9、C D 2 9 0、C D 2 9 2 / B M P R 1 A / A L K - 3、C D w 2 9 3 / B M P R 1 B / A L K - 6、C D 2 9 4、C D 2 9 5 / L E P R、C D 2 9 6、C D 2 9 7、C D 2 9 8、C D 2 9 9 / D C - S I G N R、C D 3 0 0 a、C D 3 0 0 b、C D 3 0 0 c、C D 3 0 0 e、C D 3 0 0 f、C D 3 0 1 / C L E C 1 0 A、C D 3 0 2 / C L E C 1 3 A、C D 3 0 3、C D 3 0 4 / ニューロピリン - 1、C D 3 0 5、C D 3 0 6 / L A I R 2、C D 3 0 7 a、C D 3 0 7 b、C D 3 0 7 c、C D 3 0 7 d、C D 3 0 7 e、C D 3 0 9 / V E G F R 2 / F l k - 1、C D 3 1 2、C D 3 1 4 / N K G 2 D、C D 3 1 5、C D 3 1 6、C D 3 1 6、C D 3 1 7、C D 3 1 8、C D 3 1 9 / C R A C C / S L A M 7、C D 3 2 0、C D 3 2 1 / J A M - A / F 1 1 R、C D 3 2 2 / J A M - B、C D 3 2 4 / E - カドヘリン / C D H 1、C D 3 2 5 / C D H 2 / N - カドヘリン、C D 3 2 6 / E p C A M、C D 3 2 7、C D 3 2 8、C D 3 2 9、C D 3 3 1 / F G F R 1、C D 3 3 2 / F G F R 2、C D 3 3 3 / F G F R 3、C D 3 3 4 / F G F R 4、C D 3 3 5、C D 3 3 6 / N C R 2 / N K p 4 4、C D 3 3 7 / N C R 3 / N k p 3 0、C D 3 3 8、C D 3 3 9 / J A G 1 / ジャギド 1、C D 3 4 0 / H E R 2 / E r b B 2、C D 3 4 4、C D 3 4 9、C D 3 5 0 / F r i z z l e d - 1 0 / F Z D 1 0、C D 3 5 1、C D 3 5 2、C D 3 5 3、C D 3 5 4、C D 3 5 5、C D 3 5 7、C D 3 5 8、C D 3 5 9、C D 3 6 0、C D 3 6 1、C D 3 6 2 および C D 3 6 3。

【 0 0 6 1 】

特異的細胞株を同定するためには、特異的マーカーの組合せが使用される。例えば、主要マーカー（すなわち、細胞株を代表するマーカー）は、以下の通りである：

- T細胞用のC D 3、C D 4 およびC D 8、
- B細胞用のC D 1 9 およびC D 2 0、
- 樹状細胞用のC D 1 1 c およびC D 1 2 3、
- ナチュラルキラー（NK）細胞用のC D 5 6、
- 造血幹細胞用のC D 3 4、
- 単球 / マクロファージ用のC D 1 4 およびC D 1 3 3、
- 顆粒球用のC D 6 6 b、
- 血小板用のC D 4 1、C D 6 1 およびC D 6 2、
- 赤血球用のC D 2 3 5 a、
- 内皮細胞用のC D 1 4 6、および
- 上皮細胞用のC D 3 2 6。

【 0 0 6 2 】

1つの有利な実施形態において、本発明は、検出ステップbがフローサイトメトリーによって実施される、上記で定義された方法に関する。

【 0 0 6 3 】

フローサイトメトリーは、通常は細胞である単一の粒子が流体流内を流れるにつれて、光ビームを介してその多数の物理的特性を同時に測定し、次に解析する技術である。測定される特性としては、粒子の相対的サイズ、相対的粒状性または内部複雑性、および相対的蛍光強度が含まれる。これらの特性は、細胞または粒子がいかにして入射レーザー光を散乱させ蛍光を発生するかを記録する光学 - 電子結合システムを用いて決定される。

【 0 0 6 4 】

フローサイトメータ内で、粒子は流体流内のレーザーインターセプトまで運ばれる。サイズが0.2 ~ 150マイクロメートルである懸濁した粒子または細胞は全て、解析に好適である。中実組織由来の細胞は、解析前に脱凝集されなければならない。粒子がある流体流の部分は、標本コアと呼ばれる。散乱蛍光光は、適切に位置づけられたレンズにより収集される。ビームスプリッタおよびフィルターの組合せが、散乱した蛍光光を適切な検

出器に誘導する。検出器は、それを打撃する光信号に比例する電子信号を生成する。各粒子または事象についてのリストモードデータが収集される。各事象の特性またはパラメータは、その光散乱および蛍光特性をベースとするものである。データは収集され、コンピュータ内に記憶される。このデータを解析してレーザーインターセプトを介して標本内部の亜集団についての情報を得ることができ、これらはレーザー光を散乱させる。

【0065】

光散乱は、粒子が入射レーザー光を偏向させる時に発生する。この発生程度は、粒子の物理的特性、すなわちそのサイズおよび内部複雑性に左右される。光散乱に影響を及ぼす因子は、細胞の膜、核、および細胞内部のあらゆる顆粒状物質である。細胞の形状および表面トポグラフィも同様に、全光散乱に寄与する。前方散乱光(FSC)は、細胞表面積またはサイズに比例する。FSCは、主に回折光の測度であり、フォトダイオードによる前方向の入射レーザービームの軸をわずかに外して検出される。FSCは、所定のサイズより大きい粒子をそれらの蛍光とは独立して検出する好適な方法を提供し、したがって、多くの場合、信号処理を始動させるため、免疫表現型検査において使用される。

10

【0066】

側方散乱光(SSC)は、細胞の粒状性または内部複雑性に比例する。SSCは、屈折率の変化が存在する細胞内部の任意の界面において発生する主に屈折し反射した光の測定である。SSCは、収集レンズによりレーザービームに対しおよそ90度で収集され、次にビームスプリッタにより適切な検出器に方向転換される。

20

【0067】

各々が488nmで励起され、ピーク発光波長が互いに極端に近いものでないならば、同時に2つ以上の蛍光色素を使用することができる。FITCとフィコエリスリン(PE)の組合せが、これらの基準を満たす。488nmのレーザーを用いて5つの異なる蛍光色素(FITC、PE、ECD、PC-Cy5またはPC-Cy5.5およびPC-Cy7)を検出できること、および紫色レーザーを用いて7つの異なる蛍光色素(例えば上述の7つのBrilliant Violet(登録商標)染料)を検出できることを指摘しておかなければならない。PEの吸収最大値は488nmにあるわけではないが、蛍光色素はこの波長で充分励起されて、検出用の適切な蛍光発光を提供する。より重要なことに、ピーク発光波長は、FITCについては530nm、PEについては570nmである。これらのピーク発光波長は、充分離れており、そのため各信号を別個の検出器により検出することが可能である。検出された蛍光信号の量は、粒子上の蛍光色素分子の数に比例している。

30

【0068】

フローサイトメータは、次の5つの主要構成要素を有する：

- フローセル：検知のため光ビームを通して単一のファイルを通過するように細胞を担持し整列させる液体流(シース液)
- 測定システム：一般的に使用されるのは、インピーダンス(または伝導度)測定および、光学系-ランプ(水銀、キセノン)；高出力水冷式レーザー(アルゴン、クリプトン、色素レーザー)；低出力空冷式レーザー(アルゴン(488nm)、赤色HeNe(633nm)、緑色HeNe、HeCd(UV))；光信号を結果としてもたらずダイオードレーザー(青色、緑色、赤色、紫色)である。
- 検出器およびアナログ-デジタル変換(ADC)システム：これは、FSCおよびSSCならびに光由来の蛍光信号をコンピュータによって処理できる電子信号へと変換する。
- 増幅システム：線形または対数増幅、
- 信号解析用コンピュータ。

40

【0069】

フローサイトメータを用いた標本由来のデータの収集プロセスは、「取得」と称される。取得は、フローサイトメータに対して物理的に接続されたコンピュータおよびフローサイトメータとのデジタルインターフェースを取扱うソフトウェアを介して行われる。ソフ

50

トウェアは、検査中の標本のためにパラメータ（すなわち電圧、補償など）を調整することができ、パラメータが正しく設定されることを保証するため標本データの取得中に初期標本情報を表示することも補助する。初期のフローサイトメータは、一般に実験的デバイスであったが、技術的進歩により、臨床および研究の両方のさまざまな目的で使用するための広範な応用が可能となった。これらの開発により、計装、解析ソフトウェアならびに蛍光標識された抗体などの取得において使用される試薬に対する莫大な市場が発展した。

【0070】

最新計器は通常、多数のレーザーおよび蛍光検出器を有する。市販の計器では、レーザー数4個または5個および蛍光検出器数18個が現在の記録である。レーザーおよび検出器の数を増大させることにより、多数の抗体標識が可能になり、その表現型マーカーにより標的集団をより精確に同定することができる。一部の計器はさらに個別の細胞のデジタル画像を撮影することさえでき、細胞の表面の内部またはその上の蛍光信号の場所の解析を可能にする。

10

【0071】

別の有利な実施形態において、本発明は、k平均派生アルゴリズムまたはDBSCANアルゴリズムであるクラスタリングアルゴリズムを実施することによりクラスタ化が達成される、上記で定義した方法に関する。

【0072】

上述の通り、有利なアルゴリズムはk平均、詳細には、修正k平均アルゴリズム：DenseKmeansである。

20

【0073】

本発明において、高密度集団または領域に属する細胞をクラスタ化できる任意のアルゴリズムを実行することができる。

【0074】

DenseKmeansアルゴリズムについては、実施例1に説明されている。

【0075】

さらに別の有利な実施形態において、本発明は、前記細胞大集団が、血液標本、脳脊髄液、羊水、気管支肺胞洗浄液、母乳、頸腔部液などのあらゆる動物またはヒトの体液に由来するものである、上記で定義された方法に関する。

30

【0076】

より有利には、本発明は、血液標本中の成熟内皮細胞の亜集団を同定することを目的とし、大集団の細胞が、少なくともCD45、CD105およびCD146というマーカーで標識されている、上記で定義された方法に関する。

【0077】

より有利には、本発明は、血液標本中の前駆内皮細胞の亜集団を同定することを目的とし、大集団の細胞が、少なくともCD45、CD34、CD133、およびCD309というマーカーで標識されている、上記で定義された方法に関する。

【0078】

有利には、本発明は、血液標本中の成熟または前駆内皮細胞の亜集団を同定することを目的とし、大集団の細胞が、少なくとも7AAD、CD31、CD45、CD105、CD146、CD34、CD133、CD144およびCD309というマーカーで標識されている、上記で定義された方法に関する。

40

【0079】

より有利には、本発明は、上皮細胞の亜集団を同定することを目的とし、大集団の細胞が少なくともCD326、CD45、サイトケラチンに向けられた抗体といったマーカーで標識されている、上記で定義された方法に関する。

【0080】

有利には、本発明は、上皮細胞の亜集団を同定することを目的とし、大集団の細胞が少なくともCD44、CD326、CD45、サイトケラチンに向けられた抗体といったマーカーで標識されている、上記で定義された方法に関する。

50

【0081】

より有利には、本発明は、調節B細胞の亜集団またはエプスタイン・パールウィルス（EBV）に感染した記憶B細胞の亜集団を同定することを目的とし、大集団の細胞が少なくともCD27、CD24、CD19、およびIL-10というマーカーまたはCD27、CD19およびEBV抗原に向けられた抗体というマーカーでそれぞれ標識されている、上記で定義された方法に関する。

【0082】

有利には、本発明は、調節B細胞の亜集団を同定することを目的とし、大集団の細胞が少なくともCD38、CD27、CD24、CD19、CD45、IL-10、IgDおよびCD5というマーカーで標識されている、上記で定義された方法に関する。

10

【0083】

有利には、本発明は、エプスタイン・パールウィルス（EBV）に感染した記憶B細胞の亜集団を同定することを目的とし、大集団の細胞が少なくともCD38、CD27、CD19、IgDおよびEBVに向けられたモノクローナル抗体またはポリクローナル抗体というマーカーで標識されている、上記で定義された方法に関する。

【0084】

より有利には、本発明は、調節T細胞の亜集団を同定することを目的とし、大集団の細胞が、少なくともCD4、CD25、Foxp3およびサイトカインに向けられた抗体というマーカーで標識されている、上記で定義された方法に関する。

【0085】

より有利には、本発明は、ヒト免疫不全ウイルス（HIV）に感染したCD4+T細胞の亜集団を同定することを目的とし、大集団の細胞が少なくともCD4、CD3、CD25およびHIV抗原に向けられた抗体というマーカーで標識されている、上記で定義された方法に関する。

20

【0086】

有利には、本発明は、 $\alpha = 10^{-1}$ であり、 $N = 10$ である、上記で定義された方法に関する。N = 500である場合に、いくつかの有利な結果が得られる。

【0087】

換言すると、本発明は、

a. 前記大集団の細胞をn種の試薬に曝露するステップであって、前記n種の試薬が、前記大集団の各細胞のn種の異なる構成要素の存在、不存在または量の検出を可能にするステップであって、nは2以上であるステップと；

30

b. 各細胞をn次元空間内部の特定の位置に割り当てるために、前記大集団に属する各細胞について前記n種の試薬を検出するステップと；

c. 細胞をクラスタ化によってk個の異なるクラスタにグループ化するステップであって、各クラスタが中心 C_k および半径Dにより特徴づけされており、クラスタ化は、前記大集団に属する細胞の20%~90%が前記k個のクラスタの1つに割り当てられるようなものであり、kおよび C_k のパラメータが、前記規定のクラスタに割り当てられる細胞の前記百分率に左右されるステップと；

d. 隣接するクラスタをグループ化してより大きいクラスタを得、ここで隣接するクラスタは、2つのクラスタの中心C間のユークリッド距離が半径Dの2倍よりも小さいようなものであり、かつより大きい前記クラスタの中心 C_{1k} ならびにより大きい前記クラスタに属する細胞の共分散行列を推定するステップと；

40

e. 前記n次元の各々の中より大きいクラスタの半径を係数 $\alpha = 0.1$ だけ増加させることによって各々の拡大クラスタについてのすべり領域を定義づけし、前記すべり領域に属する各細胞について、マハラノビス距離を計算するステップと；

f. $D_{1k} (1 + \alpha)$ よりも小さいマハラノビス距離を有する1セットの細胞に属する細胞の数を推定し、前記セットの密度を測定するステップであって、前記セットの細胞が、すべり領域に属するがより大きいクラスタには属していない細胞に対応しており、したがって、

50

- 密度が値 $N = 10$ より高い場合、前記セットは前記特異的細胞を含まないものとみなされ、前記セットの密度が前記値 N よりも低くなるまでステップ f および g が p 回反復され、前記セットは $D_{1k} (1 + \quad)^P$ より小さいマハラノビス距離を有する細胞により定義され、

- 前記セットの密度が N 以下である場合、前記セットは前記特異的細胞を含んでいる、ステップと；

を含む上述の方法に関する。

【0088】

本発明は同様に、癌、血管および免疫病理および感染性疾患を含めた病理のインビトロ診断またはインビトロ予後診断のための方法において、上記で定義された n 次元空間内の細胞大集団中の特異的細胞の亜集団を同定するステップを含む方法にも関する。

【0089】

有利には、本発明は、癌、血管および免疫病理および感染性疾患を含めた病理のインビトロ診断またはインビトロ予後診断のための方法において、

a. 前記大集団の細胞を n 種の試薬に曝露するステップであって、前記 n 種の試薬が、前記大集団の各細胞の n 種の異なる構成要素の存在、不存在または量の検出を可能にするステップであって、 n は 2 以上であるステップと；

b. 各細胞を n 次元空間内部の特定の位置に割り当てるために、前記大集団に属する各細胞について前記 n 種の試薬を検出するステップと；

c. 細胞をクラスタ化によって k 個の異なるクラスタにグループ化するステップであって、各クラスタが中心 C_k および半径 D により特徴づけされており、クラスタ化は、前記大集団に属する細胞の 20% ~ 90% が前記 k 個のクラスタの 1 つに割り当てられるようなものであり、 k および C_k のパラメータが、前記規定のクラスタに割り当てられる細胞の前記百分率に左右されるステップと；

d. 隣接するクラスタをグループ化してより大きいクラスタを得、ここで隣接するクラスタは、2つのクラスタの中心 C 間のユークリッド距離が半径 D の 2 倍よりも小さいようなものであり、かつより大きい前記クラスタの中心 C_{1k} ならびにより大きい前記クラスタに属する細胞の共分散行列を推定するステップと；

e. 前記 n 次元の各々の中より大きいクラスタの半径を 0.01 から 0.1 まで変動する係数 だけ増加させることによって各々の拡大クラスタについてのすべり領域を定義づけし、前記すべり領域に属する各細胞について、マハラノビス距離を計算するステップと；

f. $D_{1k} (1 + \quad)$ よりも小さいマハラノビス距離を有する 1 セットの細胞に属する細胞の数を推定し、前記セットの密度を測定するステップであって、前記セットの細胞が、すべり領域に属するがより大きいクラスタには属していない細胞に対応しており、したがって、

- 密度が値 N より高い (N は 10 から 1000 まで変動する) 場合、前記セットは前記特異的細胞を含まないものとみなされ、前記セットの密度が前記値 N よりも低くなるまでステップ f および g が p 回反復され、前記セットは $D_{1k} (1 + \quad)^P$ より小さいマハラノビス距離を有する細胞により定義され、

- 前記セットの密度が N 以下である場合、前記セットは前記特異的細胞を含んでいる、ステップと；

g. 癌、血管および免疫病理および感染性疾患を含めた病理を表わす前記特異的細胞のうちの細胞の存在を同定するステップと；

を含む方法に関する。

【0090】

上記で定義された全ての方法において、有利には、補足的ステップが実施される。ステップ g の前に発生するこのステップ f' は、高密度クラスタに属する細胞、すなわちより大きいクラスタ、有利には p 個のすべり領域により拡大されたより大きいクラスタに属する細胞とみなされている細胞を除去することからなる。

【0091】

本発明は同様に、先に定義された方法のステップc～fを実施できるようにする適切な媒体上のコンピュータプログラムにも関する。

【0092】

換言すると、本発明は、上述の実施形態のいずれか1つに記載のステップc～fの各々を実施するように適応された命令を含む、コンピュータ可読媒体上に記憶されコンピュータ上で実行されるコンピュータプログラムに関する。

【0093】

本発明の機能の一部または全てがソフトウェアを用いて実施される場合、このソフトウェア（コンピュータプログラム）は、コンピュータを用いて読取り可能な記憶媒体上に記憶された形で提供され得る。本発明については、「コンピュータ可読記録媒体」は、フロッピーディスク（登録商標）またはCD-ROMなどのポータブルフォーマットの記録媒体に限定されず、さまざまなタイプのRAMおよびROMなどのコンピュータ内の内部メモリーデバイス内またはハードディスクなどのコンピュータに固定された外部記録デバイス内にも格納され得る。

10

【0094】

本発明は同様に、

- a. 前記大集団の各細胞のn種の異なる構成要素の存在、不存在または量を検出するためのn種の試薬であって、nは2以上である、試薬と、
 - b. 上記で定義された前記n種の異なる細胞マーカーの検出用手段と、
 - c. 上記で定義されたコンピュータプログラムと、
- を含むキットにも関する。

20

【図面の簡単な説明】

【0095】

【図1A-B】人工的に生成されたデータ上でのRAREを用いた希少事象の検出を示す。

【図1A】処理前のデータセットを表わす。

【図1B】本発明に係る処理後のデータセットを表わす：2つの希少事象が存在する。すなわち1つは、疎かつ大域的であり、もう1つは高密度かつ局所的である。

30

【図2A-C】本発明に係る方法を示す。

【図2A】オリジナルデータを表わす：希少事象は、全データコレクション(X)の1%を含む。

【図2B】高密度領域のコアを除去した後のデータサブセット(X_{KEEP})を表わす。

【図2C】高密度集団の除去後の希少事象(X_{RARE})を表わす。

【図3A-I】図2からのオリジナルデータを考慮したDenseKMeans内の変動するDMAXおよびKl： $(A、B、C) DMAX = 1 : 4、Kl = 4$ ； $(D、E、F) DMAX = 1 : 2、Kl = 6$ ； $(G、H、I) DMAX = 1、Kl = 8$ 。大きい点はクラスタ中心を表わす。各ケースについての初期、中間および最終ステップは、高密度領域のコアに向かうクラスタ中心の収束を示し、外れ値に対するk平均の初期感度を削除している。

40

【図4】灰色の点は、DenseSlideを通して除去される。図3の場合と同じDMAXおよびKIの組合せが使用される。図4Cでは、8個のクラスタのうち7個だけが残され、1つは、DenseKMeansにおける密度条件(NI)を満たさないことから除去されている。

【図5】図2からの人工データセットについてのDMAXおよびKIの初期化を表わす。

【図6】人工的に生成されたデータに対しLOF（上および中間の図版）およびRARE（下の図版）を適用することによる外れ値検出例を表わす。RARE内の大きな点は、DenseKMeans由来のクラスタ中心を示す。

【図7A-B】方法の応用を表わす。

【図7A】検出チャネル対(FL1、FL2、FL6、FL7、FL8およびFL9)に

50

したがったオリジナルデータを表わす。

【図7B】検出チャンネル対 (FL1、FL2、FL6、FL7、FL8およびFL9) にしたがった希少事象を表わす。

【図8】図7由来のフローサイトメトリーデータセットについてのDMA XおよびKIの初期化を表わす。

【図9】さまざまな半径値についてのLOCIの外れ値度スコアを表わす。希少事象内の全ての点は1スコアを有し、異常として同定され得ない。

【実施例】

【0096】

実施例1：

序：

外れ値は、それが異なる機序によって生成されたという疑義を喚起するほどに他の観察事実から逸脱した観察事実である (Hawkins 1980)。同様に、希少事象 (外れ値クラスタ (Rocke 1996)、クラスタ化された異常 (Liu 2010; Liu 2012)、異常コレクション (Dai 2012)、マイクロクラスタ (Bae 2012)) は、それが異なる機序によって生成されたという疑義を喚起するほどに他の観察事実群から逸脱した観察事実群である。

【0097】

高い再現性を有する、すなわち偽陰性が全くない希少事象の検出は、希少事象を見落とすことのコストが著しく高い分野に固有のものである。最も代表的な例は、例えば血液標本中の病的細胞群の見落としに由来するコストが、健全な細胞群を病的として分類することに由来するコストよりも著しく高い、すなわち偽陰性に比べて偽陽性を優位とする医療分野である。生物学的監視中の疾病の発生 (Shmueli 2010)、クラスタ化された発病の突発 (Liu 2010) またはソーシャルメディア内のスパム送信者 / 不正なレビューアー群 (Dai 2012) は、希少事象を検出することがそれらの検出に由来するコストよりも優位である他のシナリオ例である。

【0098】

単一のまたはクラスタ化された、異常とは、正常な挙動に対して正常でないものとみなされる事象である (Chandola 2009)。いずれのタイプの異常の場合でも、解決すべき問題は、正常性を定義づけることである。単一外れ値については、他のデータインスタンスとの近傍類似性、または距離、分布の見地からみて正常性が定義づけされる。空間的異常の場合、空間を異常にするのは特定の空間領域内における異常の発生である。集合的異常の場合、個別のインスタンスは正常であるが、集合を異常にしているのは異常の同時発生である。希少事象の場合、それらを異常にしているのは、他のデータ垂集団との関係における相対的に小さい異常のサイズである。集合的異常とは異なり、希少事象内に含まれる全てのインスタンスが1つの異常である。その著しく小さいサイズを除いて、希少事象の他のデータ特性、例えばフィーチャ分布または空間的位置づけは、正常なデータ垂集団に比べていかなる判別情報も担持していない。発明者らは、図1において希少事象検出の一例を考慮している。データ分布は10,000点の2つの正常集団と2つの希少事象、すなわち正常集団とはかけ離れた10点のより疎な希少事象すなわち大域的異常と、正常集団の1つに近い20点のより高密度の希少事象すなわち局所的異常とを含む。図1(b)は、希少事象を残りのデータから単離する1つのアプローチ、RAREの出力を示している。

【0099】

外れ値およびクラスタの両方と共通の特性を共有して、希少事象の検出は、外れ値検出と強度に不均衡 / 非平衡なクラスタリングの間の境界に存在する。クラスタリングおよび外れ値検出アルゴリズムは両方共、概して、陽性例すなわち希少事象を陰性として誤って分類する傾向がある。非平衡データについてのアルゴリズムは主として、再サンプリング、原価重視または1クラス学習方法 (Chawla 2004) を用いて問題が一般に取扱われる非平衡なトレーニングデータの存在下における分類の問題向けに監視下のシナリオ

10

20

30

40

50

(Tang 2009)の中で提供されてきた。しかしながら非監視下のシナリオでは、クラスタリングアルゴリズムが概してクラスタサイズを平衡化する傾向にあることから、問題はさらに取扱い難いものとなる。例えば、k平均は、より優れた精度のためのトレードオフとしてクラスタサイズの変動を低減させる傾向をもつ(Xiong 2006)。スペクトルクラスタリングでは、RatioCutおよびNcut(Luxburg 2007)の両方が、カット値を最小化することよりもクラスタの平衡化の方により重きを置いている。両方のアルゴリズム共、導入された平衡化制約条件を通して、初期MinCut解の外れ値感度を扱うことを提案している。その一方で、外れ値/異常検出アルゴリズムは、単一の異常を発見するのに極めて有効である。文献中で、異なるアプローチ(密度ベース、距離ベース、分布ベースのもの)が提案されてきている。最も一般的な外れ値検出アルゴリズム、LOF(Breunig 2000)すなわち局所的外れ値因子は、各点の局所的密度をその近傍の点の局所的密度に対し比較することによって得られる外れ値度スコアにしたがってトップk外れ値のリストを出力する。LOFでは、結果の質は、主として近傍の構成(パラメータMinPts)により左右される。本明細書において、発明者らは、外れ値検出とクラスタリング方法の間のこのギャップに対処している。偽陰性を回避すること、すなわち真陽性を見落しを回避することが主要な課題であることから、発明者らは、密度をベースとする逆方向あるいはボトム-アップアプローチ、すなわち最も高密度の領域から最も低密度の領域に向かうアプローチを提案している。一般的な外れ値検出方法は、順方向またはトップ-ダウンアプローチを使用する。すなわち、これらの方法は、外れ値度閾値スコアにしたがってトップk外れ値を取る。本明細書は以下のように組織されている。第2節は、大きいデータセット内の希少事象を発見するための文献精査に充てられている。第3節は、RAREフレームワークを紹介している。発明者らは、まず最初に、空間の高密度領域内に存在する点のみを発見しクラスタリングするように構想されたk平均の修正変形形態であるDenseKMeansを用いてクラスタリングを実施している。第2のステップにおいて、発明者らは、密度ベースのすべり領域を用いてDenseKMeansにより発見された高密度領域を漸進的に増大させている。すべり領域内部の密度が密度条件を満たすことができなくなった時点で直ちに、発明者らは高密度領域の境界線に到達したとみなす。希少事象はこれらの境界線の外側に存在する。第4節では、合成データおよび実データの両方についての実験によって、RAREが、他の方法では失敗している希少事象を発見することができるということが示されている。発明者らは本明細書を、第5節における結論とさらなる展望についての論述で締めくくっている。

【0100】

2. 関連する研究作業

大きいデータセット内の希少事象の検出のために、文献中の異なるアプローチ(Chandola 2009; Ertöz 2003; Ester 1996; He 2003; Liu 2010; Liu 2012; Papadimitriou 2003; Zhu 2010)が提案されてきた。いくつかの技術が、クラスタベースの異常検出としてそれにアプローチしている(Chandola 2009); 正常インスタンスは大きく高密度のクラスタに属し、一方異常は小さいまたは疎なクラスタに属している。このような方法は、クラスタリングアルゴリズムの出力に依存している。CBLOF(He 2003)は第1に、任意のクラスタリング方法を用いてクラスタリングを実施し、その後、予め定義された閾値に基づいて小クラスタを大クラスタから分離する。この閾値を用いて、CBLOFは、クラスタのサイズおよび最も近いクラスタ中心までの距離の両方を考慮に入れることによって、クラスタベースの局所的外れ値因子(CBLOF)外れ値度スコアを定義づけする。全体として、このような技術の性能は、初期クラスタリングの選択および質に強く依存している。

【0101】

明示的クラスタサイズ制約条件の利用は、データセット内の希少事象の検出を取扱うのに使用可能な別の解決法(Zhu 2010)である。文献中では、クラスタの平衡化に集中する傾向にあるが、このアプローチは、異なるクラスタサイズを有するパーティショニ

ングの生成を可能にする。それは、データ内の各クラスタのサイズについての推測的な知識が予め公知である場合に、非常に有用であり得る。それでも、このような忠実な情報の恩恵を享受できるのはわずかな利用分野に過ぎない。

【0102】

第3のアプローチは、単一の外れ値検出アルゴリズムを使用するかまたは適応させ、それらを外れ値のマイクロクラスタの検出に好適なものにすることである。LOF (Breunig 2000)では、外れているクラスタの検出は、局所的近傍を定義する最も近い近傍系 $MinPts$ の数の選択によって左右される。非常に小さいクラスタの検出には、1つのクラスタ内の全ての点を格納するのに充分大きい、すなわちクラスタのサイズよりも大きい $MinPts$ が必要とされる。LOCI (Papadimitriou 2003)は、多粒度偏差因子 (MDEF) を定義し、その近傍系の近傍サイズとは著しく異なる近傍サイズを有する点として外れ値を同定する。LOFと同様に、LOCIは、近傍サイズの適切な選択に依存するが、ただし、LOFとは異なり入力パラメータとして近傍の最大半径を必要とする。

10

【0103】

別の異なる方向は、正常インスタンスがデータ内の1クラスタに属している一方で、外れ値はいかなるクラスタにも属していないと考えることにある (Chandola 2009)。このアプローチでは、全ての点を強制的にクラスタの1つに属するようにすることのない方法 (DBSCAN (Ester 1996)、SNNベースのクラスタリング (Ertoz 2003)) の使用が求められる。DBSCAN (Ester 1996) は、最も一般的な密度ベースのクラスタリングアルゴリズムである。それは、クラスタが異なるサイズおよび形状のものとなることができるようにする密度到達可能性という新規の概念に基づいてクラスタを構築する。ただし、DBSCANの性能は、クラスタが異なる密度のものであり、そのランタイム複雑性およびメモリの両方共が高い $O(n^2)$ である場合に低いものである。

20

【0104】

比較的最近のコンセプトである「単離」が、大部分の外れ値検出方法において使用される距離および密度のコンセプトに対する1つの代替案として提案された (Liu 2008; Liu 2010; Liu 2012)。単離の概念は、「わずかでかつ異なるものである」という異常の特性に依存している。このコンセプトに依存する2つの方法; $iForest$ (Liu 2008; Liu 2012) および $SCiForest$ (Liu 2010) が、トレーニング段階で、データのサブサンプリングを用いて t 個の2分木の森を構築し、評価ステップにおいて、木の根元からノードまでの経路として定義づけされる各点の経路長に基づいて異常スコアを計算する。両方の方法共、大域的にクラスタリングされた異常すなわち正常な集団からかけ離れたクラスタを発見する上で有効であるものの、 $SCiForest$ だけが、局所的にクラスタリングされた異常 (Liu 2012)、すなわち正常集団に近いクラスタを検出することができる (発明者らは図1中の実施例において、両方のタイプのクラスタリングされた異常を提示した)。しかしながら、それぞれ $O(t(q + \log +))$ および $O(qnt)$ (式中 t は $iTrees$ を構築するための抽出サンプルサイズであり、 t はトレーニング段階で構築すべき木の数である) というトレーニングおよび評価の両方の段階における高い複雑性から、 $SCiForest$ は局所的なクラスタリングされた異常の存在下でのみ適応したものになっている。

30

40

【0105】

本明細書中で発明者らが提案している RARE フレームワークは、以下のものを提案する: 1) 最初に正常/高密度領域を同定することによる希少事象の検出の逆方向アプローチ; 2) 偽陰性を回避するように構想され、したがって偽陽性を受入れる、すなわち精度よりも再現性に有利に作用するように構想されたアプローチ; 3) k 平均の変形形態 (線形、拡張可能) を使用することに起因する低い複雑性; 4) 希少事象の検出を可能にするフレームワークの2つのステップにおけるより低い接続密度駆動型アプローチ。

【0106】

50

3. RAREフレームワーク

発明者らは、この節において、大きいデータセット内の希少事象の検出のための2段階フレームワークについて説明する。N個のデータ点を有するデータセットXを所与として、発明者らは、以下に記す通り希少事象を定義づけする。

【0107】

希少事象は、サイズ N_R の点のクラスタであり、ここで N_R はデータセットの合計サイズに比べ著しく小さい($N_R \ll N$)。

【0108】

希少事象内の点の数とデータセット内の点の総数の間の比 $= N_R / N$ の形で表現した場合、上述の希少事象条件は $\ll 1$ となる。この値が非常に小さいこと、すなわち $\ll 10^{-2}$ であることから、異常事象の検出の問題は、外れ値検出と強度に不均衡なクラスタリングの間の境界に置かれることになる。

【0109】

逆方向アプローチ：例示的实施例

発明者らは、図1中の一実施例を用いてRAREの逆方向アプローチを例示している。発明者らは、2つの主要な亜集団を伴うデータセットXおよび全データセットの1%に相当する希少事象を考慮している。

【0110】

最初に、発明者らは、高密度領域のコアを同定する一方で、この段階で2つの主要な問題、すなわち拡張可能性と密度を取扱おうと考える。発明者らは、データ内の亜集団の数について推測的知識を全く有していない。拡張可能性の問題を取扱うため、発明者らは、その線形的複雑性および並列化パワーの両方のため、k平均[10]を用いてデータセットをクラスタリングすることを選択する。次に、密度問題は、高密度領域内に存在する点のみがクラスタリングされるような形でk平均を修正することによって取扱われる。発明者らはこれを、k平均の再割当て段階においてクラスタ中心の推定方法を変更することによって行なう。すなわち、クラスタ中心のまわりの最大半径のところにある点のみが、中心の再計算に使用される。半径限定型アプローチは、全ての点をクラスタの1つに強制的に帰属させることはしない。すなわち一部の点は、クラスタリングされない状態に残される。亜集団の数は予め公知ではないことから、発明者らは、初期の多数のクラスタ K_1 を使用し、多数のクラスタを用いて各集団を修正させる。図1(b)は、解析のこの第1のcステップを示している。発明者らは、この実施例において $K_1 = 6$ のクラスタ中心を使用し、DenseKMeansの出力、すなわち第1のステップ後に未クラスタリング状態に残された点XKEEPをプロットする。

【0111】

第2段階では、同じ集団に属するクラスタ(すなわちこれらは第3.3節で定義される通り隣接している)は統合されて、接続した構成要素を形成する。実施例では、3つのクラスタの各群が1つの接続構成要素を形成する。2つの構成要素はこのとき、高密度領域の境界線に達するようガウスモデルを用いて漸進的に増大させられる。これらの境界線の外側にある全てのもの(XRARE)は、希少事象とみなされる。フレームワークは、真陽性すなわち希少事象と偽陽性すなわち高密度領域の境界線の近くに存在する点または外れ値の両方を検索する(図2(c))。偽陽性は、発明者らが偽陰性を回避するために行なう妥協である。

【0112】

3.2 高密度領域クラスタリング

k平均の背後にある原理は、K個のクラスタ中心の周りにデータセットXをクラスタリングする距離ベースの目的関数の最小化に依存している。ただし、この距離ベースのアプローチは、k平均を密度関連問題に対して高感度である状態に残す。発明者らが以下で提案しているk平均の変形形態(DenseKMeans)は、オリジナルアルゴリズムに以下の2つの修正を加えることによって密度問題に対処する。

【0113】

10

20

30

40

50

【数 3】

$$\begin{aligned} \min & \sum_{k=1}^K \sum_{x_i \in C_k} \|x_i - \mu_k\|^2 \\ \text{s. t. } & |C_k| > N_I \\ & \text{dist}(x_i, CC_k) < D_{\max}, \forall x_i \in C_k \end{aligned}$$

【0114】

10

1. 初期化：各々の新しい中心が全ての他の中心から最小の DMAX 距離のところに位置づけられ、各クラスタ中心に少なくとも NI 個のデータ点が割当てられるような形で、反復的にクラスタ中心を選択する。

2. 再割当て：クラスタ中心の 1 つから最大の DMAX 距離のところにある点のみを用いてクラスタ中心を再推定し、再割当て段階中に初期 NI 閾値未満のクラスタ中心を除去する。

【0115】

DenseKMeans は、表 1 にまとめられている。クラスタ中心の 1 つから最大の DMAX 距離のところにある点のみを用いたクラスタ中心の再推定によって、半径 DMAX が外れ値までの距離よりも小さいかぎり、外れ値に対する（場合によっては希少事象に対する）k 平均の感度が削除される。その上、充分高密度でないクラスタ C_k 、($\text{card}\{C_k\} < N_I$) は、再割当て段階において廃棄される。

20

【0116】

これら 2 つの修正は、k 平均が考慮する空間の領域を高密度領域のみに限定し、クラスタ中心を反復的に高密度領域のコアに向かって移動させることを可能にする。図 3 は、異なるパラメータの組合せ DMAX 対 K、すなわち 1) DMAX = 1 : 4、K1 = 4 (図 3 (a、b、c)) ; 2) DMAX = 1 : 2、K1 = 6 (図 3 (d、e、f)) ; 3) DMAX = 1、K1 = 8 (図 3 (g、h、i)) を伴う少数の例を示す。アルゴリズムのこの第 1 の段階の出力は、オリジナルデータセットを 2 つの互いに素なサブセット $X = X_{RMV} \cup X_{KEEP}$ へと分割する。すなわち、1) X_{RMV} = 最終的クラスタ中心から最大の DMAX 距離以内に入る点、2) X_{KEEP} = 最終的クラスタ中心から最大の DMAX 距離により定義づけされる領域の外側に入る点。このアプローチを用いると、高密度領域内にある点のみがクラスタリングされる。

30

【0117】

3.3 高密度領域の増強

DenseKMeans は、クラスタ/データ亜集団の実際の数よりも著しく多い初期クラスタ数 K1 を用いて高密度領域のコアを同定する。DenseKMeans の半径限定型アプローチは、以下の通り、クラスタ隣接性特性を定義することを可能にする。

【0118】

定義 2 中心 CC_k および CC_l により定義される 2 つのクラスタは、それらの中心間のユークリッド距離が $2 \times D_{\max}$ 未満である、すなわち

40

【数 4】

$$\|CC_k, CC_l\|_2 < 2 \times D_{\max}$$

である場合、隣接している。

【0119】

DenseKMeans により発見される最終的高密度クラスタの間で、隣接するクラスタが統合されて、接続構成要素を構築し、実データ亜集団のより忠実な表現を提供する。

50

【0120】

k平均およびDenseKMeansにより使用されるもののような球形モデルでは、データの固有次元性が原次元性に等しいものとみなされる。しかしながら、実際のシナリオでは、データの固有次元性（特に局所的、すなわち1データ垂集団/クラスタ）が、原次元性と等しいものであることは稀である（LevinaおよびBickel 2005）。この課題に対処するため、我々は、データの固有次元性の取扱いにより良く適応させられたモデルを用いて球形モデルの出力を処理する。最も一般的であるのはガウスモデルである。解析の第1のステップにおいては、k平均の拡張可能性の利点のため、球形アプローチが好まれた。第1のステップにおけるガウス混合モデルの使用には、Kが予め公知でないことから、Kの全ての値についての $K(D^2 + D + 1)$ パラメータの推定が求められたはずである。たとえ、節減モデル（例えば対角線モデル）で完全ガウスモデルを置換することができたとしても、希少事象を検出しようという課題は感度が高すぎるものであり、フルモデルの使用が求められる。

10

【0121】

サブセットXRMVは、接続構成要素により定義されるコア高密度領域の共分散行列jおよび平均 μ_j の両方を迅速に推定できるようにする。これらの高密度領域は、マハラノビス距離DMおよび増加パラメータsに基づいて定義されたすべり領域SRを用いて、増強させられる。すべり領域は漸進的に高密度領域の境界線に近づき、プロセスは、密度条件が満たされる（nbPoints(SR) > NS、すなわち、すべり領域の内側の点の数は予め定義された閾値NSよりも大きい）がぎり反復される。すべり領域内部の密度がこの閾値より低下した時点で、我々は高密度領域の境界線に達したものとみなす。高密度領域増強のためのアルゴリズム、DenseSlideは表2にまとめられており、パラメータDMAXおよびK1のさまざまな組合せについてのいくつかの実施例が図4に示されている。DenseSlideのパラメータはs = 0.1、NS = 10である。アルゴリズムの出力は、陽性例のサブセットXRAREを復帰させる。

20

【0122】

4. 実験

RAREの挙動および性能は、合成データおよび実データの両方に対する実験を通して、本節で例示される。第1に、パラメータの選択についての解析および論述が、小節4.2に提示されている。我々は、さまざまなパラメータ値についてのRAREの挙動を示すために2つの人工的に生成されたデータセットを使用する。小節4.3では、医療分野由来の大規模実ケースアプリケーションに対してRAREをテストし、それを2つの他の外れ値検出方法、LOFおよびLOCIと比較する。

30

【0123】

4.1 評価

「精度」および「再現性」を用いてアルゴリズムの性能を評価する。我々の主要な課題は、真陽性を見落しを回避することにあることから、このシナリオで最も重要な評価尺度となるのは「再現性」である。高い再現性には、概して低い精度が求められる。

【0124】

【数5】

40

$$P = \frac{TP}{TP + FP} = \frac{TP}{|X_{RARE}|}$$

$$R = \frac{TP}{TP + FP} = \frac{TP}{N_{RS}}$$

（式中 $|X_{RARE}|$ はアルゴリズムによって検索されるデータ点の数であり、NRSはデータ内の陽性数、すなわち希少事象のサイズである。

【0125】

50

4.2 RARE内のパラメータ

その評価全体を通し、発明者らは、パラメータの異なる値を用いて実験し、所定の利用分野について、パラメータ値の選択が異なるデータセット全体にわたって一貫性のあるものであることを観察した。(密な関係を有するパラメータである) DenseKMeans (X_{RMV})内のデータセットのおよそ80~90%を網羅するDMAXおよび K_1 の値は、優れた最終的結果を導く。これは、希少事象がデータセット全体の10~20%という残余に比較して少ない割合に相当することから、実験の間にDenseKMeansがコア高密度領域を検出できるようにしているという事実に起因するものである。

【0126】

DenseKMeans内で1クラスに求められる最小密度 N_1 は、データセット N のサイズおよびクラス K_1 の初期数により左右され、この実験全体を通して $N_1 = N / 10 \times K_1$ に定められる。DenseSlide内のすべり領域の増大するパラメータは、 $s = 10^{-2}$ に定められ、すべり領域内で求められる点の最小数は $N_s = 10$ であった。このようにして、フレームワークの2つのステップにおける5つのパラメータのうち、発明者らには2つの自由パラメータすなわちDMAXおよび K_1 が残される。

10

【0127】

発明者らは、小節3.1内の例示的実施例の場合と同じ合成データセットを使用して、DMAXおよび K_1 の値の選択について論述する。図5は、XKEEPのサイズがこれら2つのパラメータと共にどのように推移するかを示している。先に言及した通り、発明者らは、およそ10~20%の範囲内の比 $|K_{KEEP}| / N$ を導く2つのパラメータ間の値の組合せに関心をもっている。図3で、我々は、上述の範囲内で選択されたパラメータの組合せを伴う少数の実施例を示した。概して、ほぼ同じ比 $|X_{KEEP}| / N$ について、 K_1 のより高い値と D_{MAX} のより低い値が好まれるが、それは、こうして1つのクラス内に希少事象が含まれるリスクが削減されるからである。図3内の中間ステップは、DenseKMeansの半径限定型アプローチに起因して、外れ値および外れ値/希少事象の小さい群に対する k 平均初期感度を削除することによって、クラスタ中心が高密度領域のコアに向かっていかに収束するかを十分に示している。

20

【0128】

LOFとの比較

発明者らは、図1の実施例を再度考慮する。多重パラメータ値を用いて我々のアプローチ($g-j$)およびLOF($b-f$)の両方の出力を示す。LOFが希少事象の見落とし傾向を有する一方で(a, b)、RAREはより多くの点の検索傾向を有する(j)。この挙動は、希少事象を回避するという我々の課題と合致しており、したがって、偽陰性よりも偽陽性に有利に働く。

30

【0129】

パラメータ。RAREおよびLOFは、類似のパラメータ(表3)を有する。LOFは、構成段階でMinPtsを必要とする一方、RAREは、高密度領域を発見するために2つのパラメータ D_{MAX} および K_1 を必要とする。すでに言及した通り、我々は概して、DenseKMeans内のデータのおよそ80~90%を網羅するこれら2つのパラメータの組合せを選択する。2つのパラメータを有することで、柔軟性だけでなくさらなる複雑性も加わる。DenseSlideでは、RAREは、2つのパラメータ s および N_s 、すなわち、すべり領域の成長速度および求められている最小密度を有する(s は概して 10^{-1} または 10^{-2} のいずれかに定められる)。それらの影響は、LOF内のキャッシング閾値と等価であるが、それは異なるアプローチである。すなわちLOFは、トップダウンアプローチを有し、一方RAREはボトムアップアプローチを有する。しかしながら、ボトムアップアプローチは、偽陰性の回避が優先されるシナリオにおいて、より好適である。

40

【0130】

4.3 実ケース：フローサイトメトリー

フローサイトメトリーにおいて、各細胞は、細胞マーカーに応答する蛍光レベル、すな

50

わち属性により特徴づけされる。今日、フローサイトメータは、リンパ球または単球などの任意の健常患者の体内で見い出される正常な細胞集団を表わす数千万個の細胞まで計数可能である。血液病理を呈する患者において、血液標本は同様に、異常なサインすなわち細胞マーカー蛍光レベルの異常な組合せを伴う細胞のマイクロクラスタをも含んでいる。これらの希少事象のヒト検出は、2次元空間すなわち2つのマーカーの組合せを逐次調査することによって、視覚的に実施される。このアプローチは、異常細胞集団を定義づけし複雑な多変量関係に感応するものの中で、研究所間に非常に高い相互変動性(17~44%) (Bashashati 2009)を導く。

【0131】

図7は、30個の細胞の希少事象を含む752, 987個の細胞のフローサイトメトリー標本を示す。当初FL・Log属性上でのみ可視的であった希少事象は、RAREでの処理後、他の属性上でも出現する。この実験のために、我々は、 $NI = N / 100 \times K_1$ を定めた(我々には希少事象がデータセットの合計サイズよりも著しく小さいことがわかっている)。図8は、パラメータ D_{MAX} および K_1 のさまざまな値のためのアルゴリズムの第1のステップにおいてDenseKMeansにより網羅されるデータの百分率を示している。本明細書中の残りの実験について、我々は、異なる血液標本全体にわたり比 $|K_{KEEP}| / N$ を保証する。自由パラメータ D_{MAX} および K_1 、すなわち D_{MAX} および $K_1 = 40$ を選択する。

【0132】

実験1: 変動するNR。発明者らは、第1に変動する非平衡度レベルについてのRAREの性能をテストしたいと考える。この目的で、発明者らは、定められたデータセットの合計サイズを保ち、病理期の指標である希少事象のサイズを変動させる。生物学的側面において、この実験は、健常患者の細胞抽出サンプル内に血液病理由来の成長した細胞を注入することによって実施された。注入された希少集団のサイズは{f5; 10; 20; 50; 100; 500}であった。機械の誤差に起因して、注入された細胞の数と血液標本中に発見される希少細胞集団すなわち(フローサイトメトリー内の病理シグナチャに対応する)陽性例の実際のサイズとの間に差が現われる: $NR S = \{4; 14; 17; 31; 82; 359\}$ 。全データセットは $N = NH + NR S$ の細胞を含み、ここでNHは約700:000細胞であった。DenseSlideのためのパラメータは、 $s = 10^{-1}$ および $NS = 10$ として選択された。

【0133】

表4中の結果は、特定分野の専門家により提供されたシグナチャを有することがわかった陽性例NR S(第5欄)の中で、ほぼ全ての陽性例、すなわち真陽性TP(第3欄)を発見するRAREの優れた性能を示している。RAREにより戻された偽陽性FP(第4欄)のサイズは、主として、オリジナルデータセットのサイズおよび構成によって左右される。すなわちFPは、TPが増加する一方で比較的恒常にとどまる。同様に、再現性が比較的高く、希少事象のサイズと共に精度が増加することも観察される。LOFおよびLOCIとの比較。表5では、LOF(Breunigら 2000)を用いて血液標本を解析し、異なるLOFスコア閾値についての結果を示す。解析は、希少事象を含む100Kの細胞の抽出サンプルについて実施され、NRに等しくひいては希少事象内の細胞数よりも高いパラメータMinPtsを選択した。これは、LOF内のマイクロクラスタの検出のために必要である。我々は3つの閾値、すなわち2、1.5そして希少事象内の全ての細胞を検索する最小LOFスコアに対応する第3の値を選択した。LOFおよびRAREについての同じ再現性の場合、LOFは著しく低い精度を有することが観察された(例えば、それは、NR = 20についてはデータセットのおよそ80%そして再現性を高くする場合NR = 100については60%を検索する必要がある)。閾値が、 $LOF > 1.5$ である場合、はるかに低い再現性について、精度はRAREよりつねに低いものである。我々は、NR = f5; 10; 20; 50; 100; 500gの各々で標本上にLOCIを適用した(Papadimitriouら 2003)。我々は、LOCI内でさまざまな値の最大半径{3000; 4000; 5000; 6000g}を使用した。希少事象

10

20

30

40

50

内の点について毎回1スコアを得た(図9に、LOCIスコア頻度の少数の例が提示されている)。LOFおよびLOCIの両方における1スコアは、インライアを表わし、希少事象は検出され得ない。6000超の半径値については、メモリーの問題に遭遇する。

【0134】

実験2：癌および頭蓋内動脈瘤。発明者らは、実際の患者の血液標本についてそのフレームワークをテストした(癌について4件、頭蓋内動脈瘤について6件)。標本は、生物学的ベンチマーク(細胞数200~500万個)よりも著しく大きいものであった。DenseKMeansのパラメータは、生物学的ベンチマークの場合と同じであった。すなわち、異なる血液データ標本全体にわたり比 $|K_{KEEP}|/N$ が約10~20%であることを保証するために、 $DMAX=8000$ 、 $Kl=40$ 。DenseSlideにおいて、発明者らは、実データの感度を説明するため、すべり領域のより小さな増加パラメータ $s=10^{-2}$ 、すなわち高密度領域の境界線のより緩慢な接近を選択した。高い再現性 TP/N_{RS} には、低い精度 $TP/|X_{RARE}|$ が求められる。それでも、比 $|X_{RARE}|/N$ は非常に小さく(およそ $10^{-2} \sim 10^{-3}$)、これが、高い再現性での希少事象の非常に優れた単離を保証する。再現性が低いために、我々は、停止(カッティング)パラメータ $NS=\{50; 100; 500\}$ を増大させている。NSの増加に伴って、DenseSlideはより早く停止し、こうして再現性の増大および精度の減少が導かれる。結果は、癌の標本の場合に比べて、頭蓋内動脈瘤の血液標本の場合、希少事象がより容易に単離されたことを示している。

10

20

【0135】

実験3：DBSCANおよびLOFとの比較

3つの方法が求めるパラメータの比較は、表8に提示されている。LOFは、構成段階において1つのパラメータ-MinPtsしか必要としないが、一方DBSCANおよびRAREは両方共2つのパラメータを必要とし、こうしてモデルに対しより大きな柔軟性が与えられるだけでなくより大きな複雑性も加わる。RAREおよびLOFは共に、停止基準を求めるが、一方DBSCANは、クラスタリングされていない状態に残された全ての点をノイズとみなす。DBSCANでは、希少事象が多くの場合(次の実験において示されている通り)ノイズカテゴリ内に入る。RAREは2つのパラメータ、およびNS、すなわちすべり領域の成長速度と最小密度(は概して 10^{-1} または 10^{-2} のいずれかに定められる)を使用して、停止基準を定義する。それらの影響は、LOFにおけるカッティング閾値と等価であるが、それは異なるアプローチである。すなわちLOFはトップダウンアプローチであるのに対してRAREはボトムアップアプローチである。ボトムアップアプローチは、偽陰性の回避が優先されるシナリオにおいて好まれる。

30

40

【0136】

表9において、発明者らは、3つの方法のためのさまざまなパラメータ値を用いて、中間希少事象で第2の実験(752987個の標本と31の陽性例)から無作為に選択したデータ標本を解析した。発明者らは、アルゴリズムにより検索された真陽性(TP)および偽陽性(FP)の数を計算している。RAREおよびDBSCANは共に高い再現性を有する(概して100%)が、一方RAREはDBSCANよりも著しく高い精度を有する。DBSCANでは、大部分のパラメータ値について、希少事象は未クラスタリング状態に残され、希少事象の一部が別個に小さいクラスタ内でクラスタリングする2つのケース(14個および25個の点)を除いてノイズ5として分類されたサブセットに属している。DBSCANが、比較的優れた性能を得るためにMinPtsパラメータが希少事象のサイズよりも小さくなることを求める一方で、それとは反対にLOFは希少事象のサイズよりも高いMinPtsパラメータを求める。すなわちこれはLOFにおけるマイクロクラスタの検出に必要である。DBSCANはいかなる停止基準も必要としないが、LOFでは、発明者らは、カッティング閾値または外れ値数のいずれかを選択する必要がある。発明者らはここで、LOF内の各MinPts値について2つのカッティング閾値を使用し、各ケースにおいて偽陽性の数を標示する。この2つの値は、希少事象の圧倒的多数が2つの値を境とする範囲内にLOF外れ値度スコアを有するような形で選択されたも

50

のである。

【0137】

5. 結論

発明者らはここで、大きいデータセット内の希少事象を単離するために逆方向アプローチフレームワークを提案した。これらの事象のサイズは、クラスタリングおよび外れ値検出の両方のアルゴリズムによるそれらの検出を困難にしているが、それはその両方が真陽性を偽陰性として誤って分類する傾向をもつからである。RAREフレームワークは、再現性が精度より優勢である利用分野、例えば医学、ソーシャルネットワーク内の創発的役割を標的にしている。このタイプの問題における拡張可能性および密度問題を取扱うためにk平均の新たな変形形態が提案されており、すべり領域は偽陰性を回避するように構想される。発明者らは、主要パラメータDMAXおよびKLが、第1のステップで10~20%の網羅率を保证するように選択され得、より小さいDMAXおよびより大きいKLに選好性が付与されるということを示した。複雑性は、線形があって並列化により改善可能なDenseKMeansの複雑性が主として優位に立っている。

10

【0138】

実施例の文献リスト

[1] D. H. Bae, S. Jeong, S. W. KimおよびM. Lee. 「中心度および中心近接性を用いた外れ値検出」、CIKM議事録、2012中。

[2] A. BashashatiおよびR. Brinkman. 「フローサイトメトリーデータ解析方法の調査」、Advances in Bioinformatics、2009中。

20

[3] M. Breunig, H. P. Kriegel, R. T. NgおよびJ. Sander. 「LOF: 密度ベースの局所的な外れ値の同定」、ACM SIGMOD議事録、2000中。

[4] V. Chandola, A. BanerjeeおよびV. Kumar. 「異常検出: 調査」、ACM Computing Surveys、41、2009。

[5] L. Ertöz, M. SteinbachおよびV. Kumar. 「高ノイズ高次元データ内で異なるサイズ、形状および密度のクラスタを発見する」、SDM、2003。

[6] M. Ester, H. P. Kriegel, J. SanderおよびX. Xu. 「ノイズを伴う大空間データベース内でクラスタを発見するための密度ベースのアルゴリズム」、ACM SIGKDD議事録、1996中。

30

[7] Z. He, X. XuおよびS. Deng. 「クラスタベースの局所的な外れ値の発見」、Pattern Recognition Letters 24、2003。

[8] E. LevinaおよびP. J. Bickel. 「固有次元の最大尤度推定」、Advances in Neural Information Processing Systems、17、2005。

[9] U. von Luxburg. 「スペクトルクラスタリングについてのチュートリアル」、Statistics and Computing、17、2007。

[10] J. B. MacQueen. 「多変量観測の分類および解析用のいくつかの方法」、数理統計学および確率に関する第5回パークレイシンポジウム議事録、1967中。

40

[11] S. Papadimitriou, H. Kitagawa, P. GribbonsおよびC. Faloutsos. 「LOCI: 局所相関積分を用いた高速外れ値検出」、ICDE議事録、2003中。

[12] Y. Tang, Y. Q. Zhang, N. W. ChawlaおよびS. Krauss. 「極めて不均衡な分類用のSVM」、システム、人間、サイバネティクスに関するIEEE会議議事録、39: 281-288、2009。

[13] H. Xiong, J. WuおよびJ. Chen. 「k平均クラスタリングと妥当性確認評価基準の関係: データ分布の展望」、KDD、2006。

50

[1 4] S . Zhu、D . Wang および T . Li . 「サイズ制約条件を伴うデータクラスタリング」、Knowledge - Based Systems , Elsevier , 23 : 883 - 889 , 2010 .

【 0 1 3 9 】

【 表 1 】

<p>入力: $X = \{x_i\}, i = 1..N, x_i \in \mathbb{R}^D$ K_I - 初期クラスタ数 N_I - 最小点数(密度) D_{MAX} - 半径</p> <p>出力: $CC = \{CC_k\}, k = 1..K_F$ - 最終クラスタ中心 X_{KEEP} - 未クラスタリング状態に残された点のサブセット X_{RMV} - クラスタリングされた点のサブセット</p> <hr/> <p>初期化: 1': 互いに D_{MAX} 超遠くなるように反復的にクラスタ中心 CC を 選択する:</p> $\ CC_k, CC_l\ _2 > D_{MAX}, \forall k, l = 1..K_I$ <p>2': 密度条件すなわち $card\{C_k\} > N_I$ をチェックする。 3': 収束に至るまでステップ 1' と 2' を繰り返す: 全ての K_I 中心に少なくとも N_I 個の点が割当てられる。</p> <hr/> <p>DenseKMeans: 1'': 全ての中心から D_{MAX} 超遠い全ての点 X_{KEEP} を選択する。</p> $\min(x_i, CC_k) > D_{MAX}$ <p>2'': $X_{RMV} = X \setminus X_{KEEP}$ を用いてクラスタ中心を角度推定する。 3'': クラスタ中心が初期 N_I 閾値の下に入る場合 $\{card\{C_k\} < N_I\}$ それを除去する。 4'': 収束に至るまでステップ 1''-3'' を繰り返す。最大反復数に到達し、 中心は著しく変化しない。</p>	<p>10</p> <p>20</p> <p>20</p>
--	-------------------------------

表 1: DenseKMeans.

【 0 1 4 0 】

【表 2】

入力: X_{KEEP}, X_{RMV}, CC - DenseKMeans の出力
 ϵ_S - すべり領域のための増加パラメータ
 N_S - すべり領域内の点の数

出力: X_{RARE} - 希少事象

接続構成要素:
1^o: クラスタ隣接性特性を用いてグラフ $G = (CC, E)$ を構築する。
2^o: G 中の接続構成要素 G_j を発見する。
3^o: X_{RMV} を用いて、 $N(\mu_j, \Sigma_j)$ として G_j をモデリングする。

すべり領域:
1^o: $X_{RARE} = X_{KEEP}$ を初期化する。
2^o: 各 G_j についてマハラノビス距離を計算する:

$$D_M^j = \sqrt{(X_{RARE} - \mu_j)^T \Sigma_j^{-1} (X_{RARE} - \mu_j)}$$

3^o: X_{RMV} から最も遠い点よりも構成要素の中心の1つに近い点を X_{RARE} から除去する: $D_M^j(x_i) > D_{max}^j$
4^o: 各構成要素 $N(\mu_j, \Sigma_j)$ の周囲に、移動するすべり領域 $S_R(D_{max}^j, \epsilon_S)$ を創出する。
5^o: S_R 内部で X_{RARE} から点を除去する。
6^o: 密度条件すなわち $nbPoints(S_R) > N_S$ が遵守されているかぎり、ステップ 4^o および 5^o を繰返す。

10

表2: DenseSlide.

20

【 0 1 4 1 】

【表 3】

RARE	$D_{MAX} \& K_I$ (DenseKMeans)	$\epsilon_S \& N_S$ (DenseSlide)	ボトムアップ(逆方向)
LOF	$MinPts$ (近傍)	閾値またはトップ k	トップダウン(順方向)

表3 RARE対LOFにおけるパラメータ

30

【 0 1 4 2 】

【表4】

N_R	N	$\frac{ X_{KEEP} }{N}$ (%)	$ X_{RARE} $	TP	FP	P	R	N_{RS}
0	151,388	7.7%	64	5	59	7.8%	100%	5
5	646,149	8.1%	42	4	38	9.5%	100%	4
10	780,988	7.6%	54	13	39	24%	92.8%	14
20	757,234	7.5%	70	17	53	24.2%	100%	17
50	752,987	7.4%	65	30	35	46.1%	96.7%	31
100	760,842	7.2%	132	80	52	60.6%	97.5%	82
500	718,743	7.7%	415	358	57	86.2%	99.7%	359
0	696,465	10.9%	102	14	88	13.7%	100%	14
5	731,576	11.0%	98	9	89	9.1%	75%	12
10	720,945	9.9%	114	14	100	12.2%	100%	14
20	484,285	10.5%	129	25	104	19.3%	96.1%	26
50	630,341	10.4%	40	35	5	87.5%	97.2%	36
100	676,745	10.2%	142	69	77	48.5%	98.5%	70
500	516,981	11.2%	541	366	175	67.6%	98.6%	371
0	671,582	10.1%	94	8	86	8.5%	100%	8
5	707,535	10.8%	100	7	93	7%	100%	7
10	714,081	10.2%	135	13	122	9.6%	100%	13
20	621,155	11.8%	155	11	144	7%	100%	11
50	599,851	10.2%	144	26	118	18%	100%	26
100	711,801	10.5%	204	84	120	41.1%	100%	84
500	993,671	10.7%	552	312	240	56.5%	100%	312
0	737,997	12.1%	253	9	244	3.5%	90%	10
5	711,130	10.5%	118	10	108	8.4%	100%	10
10	707,199	10.3%	113	11	102	9.7%	100%	11
20	702,362	10.4%	104	16	88	15.3%	100%	16
50	620,829	10.2%	159	29	130	18.2%	100%	29
100	674,316	10.2%	165	70	95	42.4%	100%	70
500	658,590	10.1%	593	336	257	56.6%	99.7%	337
0	602,814	9.9%	131	12	119	9.1%	100%	12
5	618,192	10.5%	93	9	84	9.6%	90%	10
10	703,027	9.5%	122	13	109	10.6%	100%	13
20	701,580	9.8%	112	16	96	14.2%	94.1%	17
50	381,654	11.0%	111	25	86	22.5%	100%	25
100	719,439	11.5%	149	64	85	42.9%	98.4%	65
500	648,391	10.1%	520	317	203	60.9%	100%	317

表4 変動する $N_R = \{5, 10, 20, 50, 100, 500\}$ の各々についての5つの標本に対するRARE

【 0 1 4 3 】

10

20

30

【表 5】

N_R (MinPts)	LOFscores	$N_{retrieved}$	TP	P	R	N_{RS}
5	>2	39	0	0%	0%	4
	>1.5	1170	1	$8.5 * 10^{-2}\%$	25%	
	>1.24	7,801	4	$5.1 * 10^{-2}\%$	100%	
10	>2	33	0	0%	0%	14
	>1.5	665	4	0.6%	28.5%	
	>1.09	25,670	14	$5.4 * 10^{-2}\%$	100%	
20	>2	35	0	0%	0%	17
	>1.5	571	1	$1.7 * 10^{-1}\%$	5.8%	
	>1.00	80,292	17	$2.1 * 10^{-2}\%$	100%	
50	>2	49	1	2%	3.2%	31
	>1.5	697	3	0.4%	9.6%	
	>1.27	5,985	31	0.5%	100%	
100	>2	45	0	0%	0%	82
	>1.5	821	2	$2.4 * 10^{-1}\%$	2.4%	
	>1.03	58,268	82	$1.3 * 10^{-1}\%$	100%	
500	>2	95	0	0%	0%	359
	>1.5	2,180	0	0%	0%	
	>1.09	38,840	359	$9.2 * 10^{-1}\%$	100%	

表5 変動する N_R についての LOF

【 0 1 4 4 】

【表 6】

Pat.	N	$\frac{ X_{KEEP} }{N}$ (%)	N_S	$ X_{RARE} $	TP	FP	P	R	N_{RS}
C1	2,470,042	20.4%	50	7,117	2	7,115	$2.8 * 10^{-2}\%$	28.5%	7
			100	14,113	4	14,109	$2.8 * 10^{-2}\%$	57.1%	
			500	91,337	7	91,330	$7.6 * 10^{-3}\%$	100%	
C2	3,413,325	36.7%	50	10,787	24	10,763	$2.2 * 10^{-1}\%$	96%	25
			100	14,311	24	14,287	$1.6 * 10^{-1}\%$	96%	
			500	28,130	25	28,105	$8.8 * 10^{-2}\%$	100%	
C3	5,989,247	16.2%	50	6,654	1	6,653	$1.5 * 10^{-2}\%$	2.4%	41
			100	36,984	31	36,953	$8.3 * 10^{-2}\%$	75.6%	
			500	95,403	41	95,362	$4.3 * 10^{-2}\%$	100%	
C4	5,959,464	15.6%	50	3,071	24	3,047	$7.8 * 10^{-1}\%$	54.5%	44
			100	4,241	25	4,216	$5.9 * 10^{-1}\%$	56.8%	
			500	19,479	32	19,447	$1.6 * 10^{-1}\%$	72.7%	

表6 癌

【 0 1 4 5 】

【表 7】

Pat.	N	$\frac{ X_{KEEP} }{N}$ (%)	N_S	$ X_{RARE} $	TP	FP	P	R	N_{RS}
A1	2,524,916	22.1%	50	6,727	15	6,712	$2.2 * 10^{-1}\%$	100%	15
			100	8,987	15	8,972	$1.6 * 10^{-1}\%$	100%	
A2	4,130,539	23.2%	50	3,615	6	3,609	$1.6 * 10^{-1}\%$	75%	8
			100	5137	6	5131	$1.1 * 10^{-1}\%$	75%	
A3	4,595,598	18.5%	50	3,986	12	3,974	$3 * 10^{-1}\%$	70.5%	17
			100	6,252	16	6,236	$2.5 * 10^{-1}\%$	94.1%	
A4	1,895,261	15.7%	50	6,971	23	6,948	$3.2 * 10^{-1}\%$	92%	25
			100	13,397	23	13,374	$1.7 * 10^{-1}\%$	92%	
A5	1,899,278	15%	50	4,698	21	4,677	$4.4 * 10^{-1}\%$	100%	21
			100	7,030	21	7,009	$2.9 * 10^{-1}\%$	100%	
A6	3,039,332	17.7%	50	6,244	18	6,226	$2.8 * 10^{-1}\%$	100%	18
			100	10,906	18	10,888	$1.6 * 10^{-1}\%$	100%	

表7 頭蓋内動脈瘤

【 0 1 4 6 】

【表 8】

方法	モデルパラメータ	停止基準	アプローチ
RARE	(D _{MAX} , K _I)	(ϵ S, NS)	ボトムアップ(逆方向)
DBSCAN	(ϵ , MinPts)	-	ボトムアップ
LOF	MinPts	閾値またはトップ-k	トップダウン(順方向)

表8. RARE、DBSCANおよびLOFにおけるパラメータ

【 0 1 4 7 】

10

【表 9】

方法	パラメータ	TP	FP
RARE($D_{MAX}, K_I, \epsilon_S, N_S$)	(6000, 80, 0.1, 10)	31	193
	(6000, 100, 0.1, 10)	31	48
	(7000, 40, 0.1, 10)	31	43
	(7000, 60, 0.1, 10)	31	60
	(7000, 80, 0.1, 10)	31	57
	(7000, 100, 0.1, 10)	30	40
	(8000, 20, 0.1, 10)	31	184
	(8000, 40, 0.1, 10)	31	60
	(8000, 60, 0.1, 10)	31	22
	(9000, 10, 0.1, 10)	31	284
	(9000, 30, 0.1, 10)	31	48
	(9000, 50, 0.1, 10)	31	35
	(10000, 10, 0.1, 10)	31	51
	(10000, 30, 0.1, 10)	31	35
DBSCAN($\epsilon, MinPts$)	(5000, 10)	31	1286
	(5000, 20)	31	1998
	(5000, 30)	31	2703
	(6000, 10)	31	457(14)
	(6000, 20)	31	699
	(6000, 30)	31	934
	(7000, 10)	31	197(25)
	(7000, 20)	31	331
(7000, 30)	31	396	
LOF($MinPts, 閾値$)	(30, 1)	31	589039
	(30, 1.1)	3	132890
	(50, 1.5)	31	2133
	(50, 1.6)	8	945
	(100, 2)	31	230
	(100, 2.5)	3	54
	(150, 2.1)	31	206
	(150, 2.7)	3	43

表9 RARE、DBSCANおよびLOFの間の比較。第2欄中のパラメータ値は、第1欄からの各方法のそれぞれのパラメータに対応する。

20

30

40

50

【0148】

実施例 2 .

実験データ

材料と方法

患者および血液標本の収集

健常ドナー、頭蓋内動脈瘤を患う患者および結腸直腸癌を患う患者から EDTA の入った試験管を用いて末梢血標本を収集した。

【0149】

フローサイトメトリー

細胞解析

CEC 解析のために、血液標本を低張性溶解洗浄手順で前処理した。4 mL の血液を 50 mL 入り試験管に移し、塩化アンモニウム 0.15 M 溶液を赤血球溶解のために 1 V / 5 V で添加した。4 で 5 分の後、懸濁液を 4 で 5 分間 $400 \times g$ で遠心分離し、上清を除去し、ペレットを 20 mL の NH_4Cl 溶液で洗浄し、懸濁液を直ちに遠心分離した ($400 \times g$ 、4 、5 分間)。

10

【0150】

ペレットを RPMI 1640 溶液で洗浄し、遠心分離の後細胞を暗所で、室温 (RT) で 15 分間、以下のモノクローナル抗体混合物と共にインキュベートした：パシフィックブルー (PB) 共役型 CD31 (クローン 5.6E; Beckman - Coulter、USA) $5 \mu\text{L}$; クロームオレンジ (KO) 共役型 CD45 (クローン J.33; Beckman - Coulter) $10 \mu\text{L}$; フルオレセインイソチオシアネート (FITC) 共役型 CD34 (クローン 581; Beckman - Coulter) $5 \mu\text{L}$; フィコエリスリン (PE) 共役型 CD105 (クローン 43A4E1; Miltenyi Biotec GmbH、ドイツ) $5 \mu\text{L}$; 7 - アミノアクテノマイシン D (Beckman - Coulter) $10 \mu\text{L}$; フィコエリスリンシアニン - 7 (PC7) 共役型 CD309 (クローン KDR - 1; Beckman - Coulter) $5 \mu\text{L}$; およびアロフィコシアニン (APC) 共役型 CD146 (クローン 541 - 10B2、Miltenyi Biotec) $5 \mu\text{L}$ 。

20

【0151】

インキュベーション後、細胞を、2% の FCS で補足されたリン酸緩衝液で洗浄した ($400 \times g$ 、5 分、4)。上清を除去した後、細胞を 1 mL の PBS 中に再懸濁させた。

30

【0152】

Summit 6.1 ソフトウェアを伴う CyAn フローサイトメータ上で標本を取得した (Beckman - Coulter)。

【0153】

HUVEC の培養

セルソーターによる HUVEC の添加によって RARE の感度および再現性を解析するためにヒト臍帯静脈内皮細胞 (HUVEC) を使用した。75 cm^2 入りフラスコ内で、内皮細胞成長培地 (Promocell) 27 mL 中、37 、5% の CO_2 で、1 cm^2 あたり細胞数 5000 ~ 10000 個の平板固定密度で、HUVEC - c (Promocell GmbH、ドイツ) を培養する。

40

【0154】

それらがひとたび 70 ~ 90% の集密度に達した時点で、7.5 mL の HEPES BSS 溶液 (Promocell) を容器表面に添加して、細胞を洗浄した。HEPES BSS を吸引し、7.5 mL のトリプシン / EDTA 溶液を 2 分間添加して、HUVEC を剥離させ、7.5 mL の FCS を添加してトリプシンを中和した。懸濁液を 5 分間 $220 \times g$ で遠心分離し、90 μL の PBS 中でペレットを再懸濁させた。

【0155】

HUVEC の選別

50

H U V E C を、C E C 検出の場合と同じ混合物を用いて染色し、M o f L o A s t r i o s セルソーター (B e c k m a n - C o u l t e r) 上で取得した。ダブレットを除去した後、選別を2つのパラメータ、すなわち F C S と S S C に基づいて行った。同じモノクローナル抗体で染色された 10^6 個の末梢血単核細胞の入った標準的な 5 mL 入りフローサイトメトリー試験管内で、0、5、10、20、50、100または500個の H U V E C を分布させた。

【0156】

P B M C および H U V E C を含む懸濁液を C y A n フローサイトメータ上で解析した。

【0157】

C E C / H U V E C のシグナチャ

多数のドットプロットを用いて問題の細胞マーカー全てを解析し、それらの閾値を決定する。このゲーティング戦略によって、発明者らは問題の集団についてのシグナチャを定義づけることができる。

10

【0158】

R A R E アプローチにおいてこのシグナチャを適用するためには、このシグナチャを情報処理言語に変換しなければならない。

【0159】

X および Y の値は、データがサイトメトリソフトウェアで解析される場合に1本のチャネルに対応するが、情報処理言語では、1023チャネルは65532個の値の中に分布している。したがって、R A R E についてのシグナチャを生成するためには、発明者らは以下の等式を適用しなければならない：すなわち、

20

$R A R E \text{ 閾値} = (\text{ソフトウェア閾値} \times 65532) / 1023$

【0160】

問題の集団を同定するための R A R E シグナチャが生成される。

【0161】

結果

フローサイトメトリーによる C E C 測定の感度および再現性を測定するため、M o f l o A s t r i o s 上で0、5、10、20、50、100および500個の H U V E C を選別し、同じマルチカラーパネルで染色した単一の採血に由来する 10^6 個の末梢血単核細胞と混合した。これらの列挙は、q u i n t u p l i c a t e で実施した。C y A n フローサイトメータ上で細胞を解析し、データ解析を K a l u z a (登録商標)ソフトウェア (B e c k m a n - C o u l t e r) および R A R E フレームワーク上で実施した。結果を下表で詳述する。

30

【0162】

【表 1 0】

		計数済み細胞									
		Kaluz TM					RAREフレームワーク				
		1	2	3	4	5	1	2	3	4	5
選別済み細胞	0	3	0	0	0	0	5	14	8	9	12
	5	3	3	2	4	5	4	9	7	10	9
	10	9	6	6	4	8	13	14	13	11	13
	20	16	12	10	12	10	17	25	11	16	16
	50	25	29	16	21	22	30	35	26	29	25
	100	68	50	67	56	56	80	69	84	70	64
	500	340	303	256	305	295	358	366	312	336	317

10

【 0 1 6 3 】

発明者らは、選別されたHUVECの数とフローサイトメトリー ($R^2 = 0.9991$) およびRAREフレームワーク ($R^2 = 0.9987$) で回収したHUVECの数の間に類似の相関関係を発見したが、全てのHUVECが検出されたわけではない。これは、サイトメータの電子的中止およびセルソーターの選別の中止に起因するかもしれない。

20

【 0 1 6 4 】

臨床的有用性を有するためには、使用中の方法は、低い変動性、そして異なる病理（癌および頭蓋内動脈瘤の治療）について検定によりCECと指定された細胞が真に内皮に由来することの確証を有していなければならない。

【 0 1 6 5 】

【表 1 1】

癌	計数済みCEC	
	Kaluz TM	RARE
C1	2	7
C2	15	25
C3	38	41
C4	35	32

30

【 0 1 6 6 】

【表 1 2】

IA	計数済みCEC	
	Kaluz TM	RARE
A1	10	15
A2	3	6
A3	15	16
A4	14	23
A5	17	21
A6	24	18

40

50

【 0 1 6 7 】

計数済み C E C について、サイトメトリソフトウェア (K a l u z a (登録商標)) を用いた大量データセットについての解析と、同じデータセット内の希少事象を単離するための逆方向アプローチのフレームワーク後の間で、著しい相違はなかった ($P > 0.05$)。

【 0 1 6 8 】

参考文献

- 1 . 分散媒中の粒子を測定するための測定器用の通気形電離箱。 D i t t r i c h および G o h d e .
- 2 . 肺癌におけるフローサイトメトリーによる DNA 含量解析の診断上および生物学的意味合い。 C a n c e r R e s . , 1 9 8 3 , 4 3 , 5 0 2 6 ~ 5 0 3 2 . B u n n P ら
- 3 . 組織細胞中の抗原の位置特定。 I I . 蛍光抗体法を用いた抗原検出方法における改良。 J . E x p . M e d . 1 9 5 0 , 9 1 , 1 ~ 1 0 . C o o n s A . H . ら
- 4 . 二重標識および蛍光デジタル画像処理顕微鏡法を用いた多変量染色体解析および完全染色体解析。 C y t o m e t r y , 1 9 9 0 , 1 1 , 8 0 ~ 9 3 . A r n d t - J o v i n D . J . ら
- 5 . ヒト乳癌解析における画像サイトメトリ DNA 解析が、フローサイトメトリーによる低 S 期分画を伴う二倍体ケースにおいて予後情報を追加する可能性がある。 C y t o m e t r y , 1 9 9 2 , 1 3 , 5 7 7 ~ 5 8 5 . B a l d e t o r p ら
- 6 . 中実新生物におけるフローサイトメトリーによる 1 本鎖および 2 本鎖 RNA 測定。 C y t o m e t r y , 1 9 9 1 , 1 2 , 3 3 0 ~ 3 3 5 . E l - N a g g a r A . K . ら
- 7 . 細胞質 pH のフローサイトメトリー測定：利用可能な蛍光色素の決定的評価。 C y t o m e t r y , 1 9 8 6 , 7 , 3 4 7 ~ 3 5 5 . M u s g r o v e E . ら
- 8 . フローサイトメトリーおよび I n d o - 1 A M を使用した、多形核白血球中のサイトゾルイオン化カルシウム変動の解析。 C y t o m e t r y , 1 9 8 9 , 1 0 , 1 6 5 ~ 1 7 3 . L o p e z M . ら

10

20

【 図 1 A 】

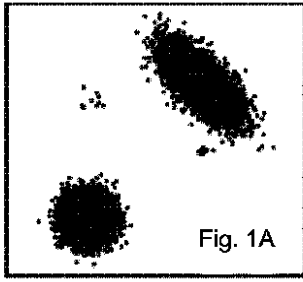


Fig. 1A

【 図 1 B 】

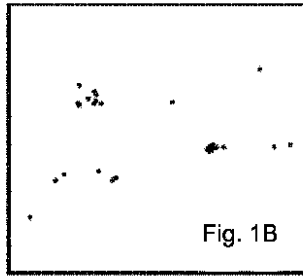


Fig. 1B

【 図 2 A 】

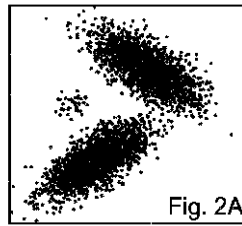


Fig. 2A

【 図 2 B 】

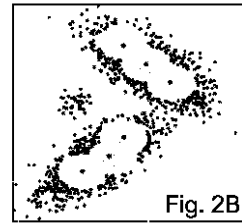


Fig. 2B

【 図 2 C 】

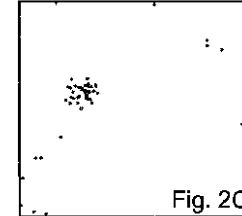


Fig. 2C

【 図 3 A 】

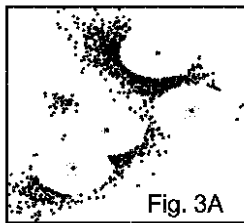


Fig. 3A

【 図 3 B 】

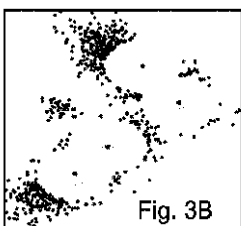


Fig. 3B

【 図 3 C 】

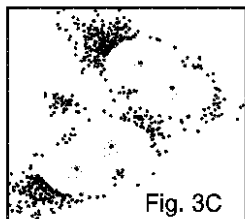


Fig. 3C

【 図 3 D 】

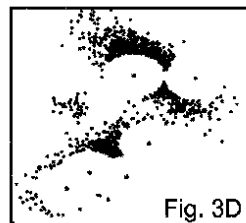


Fig. 3D

【 図 3 E 】

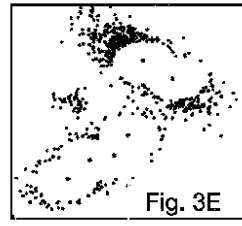


Fig. 3E

【 図 3 F 】

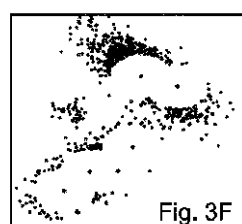
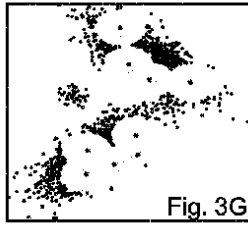
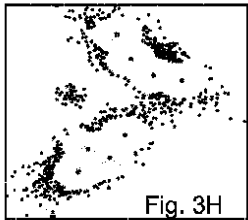


Fig. 3F

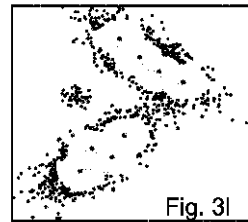
【 図 3 G 】



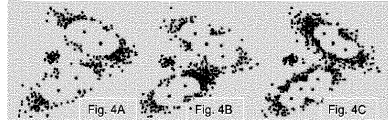
【 図 3 H 】



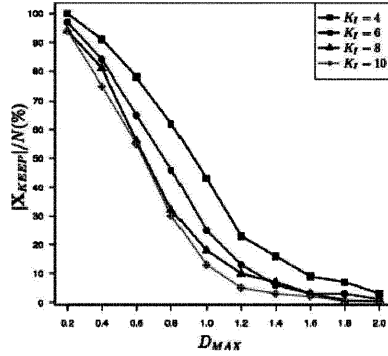
【 図 3 I 】



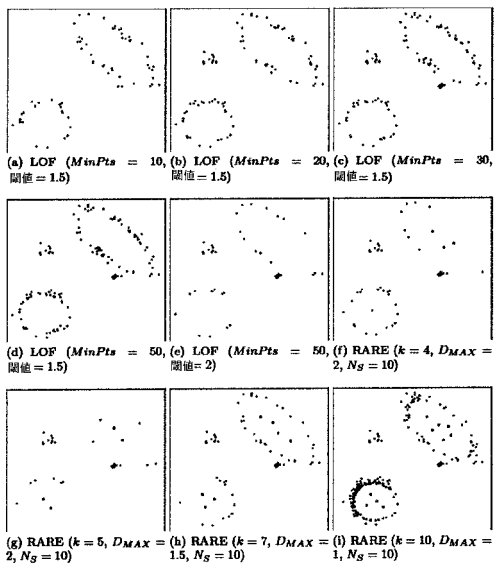
【 図 4 A - 4 C 】



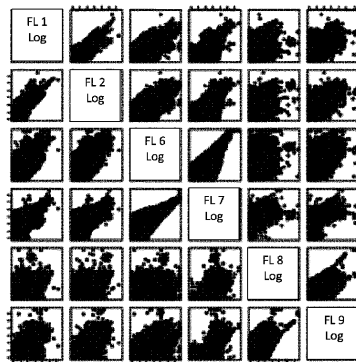
【 図 5 】



【 図 6 】



【 図 7 A 】



【 図 7 B 】

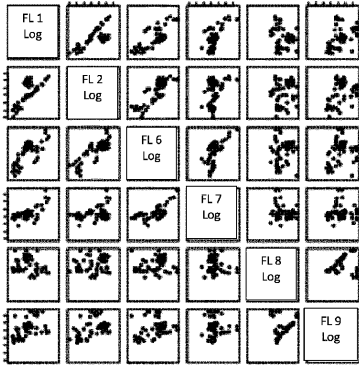


Figure 7B

【 図 8 】

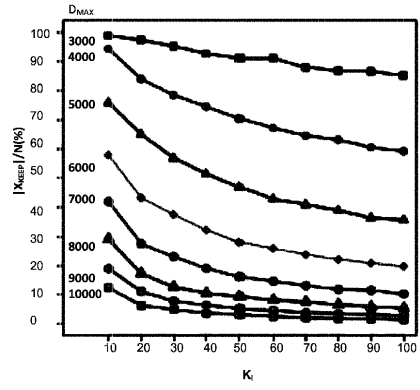


Figure 8

【 図 9 】

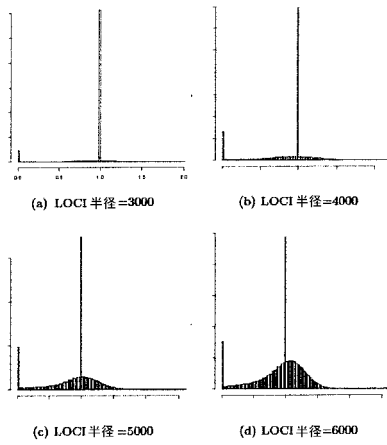


Figure 9

【 国際調査報告 】

INTERNATIONAL SEARCH REPORT

International application No PCT/EP2014/051963

A. CLASSIFICATION OF SUBJECT MATTER INV. G06F17/18 G01N15/14 ADD.		
According to International Patent Classification (IPC) or to both national classification and IPC		
B. FIELDS SEARCHED Minimum documentation searched (classification system followed by classification symbols) G06F G01N		
Documentation searched other than minimum documentation to the extent that such documents are included in the fields searched		
Electronic data base consulted during the international search (name of data base and, where practicable, search terms used) EPO-Internal		
C. DOCUMENTS CONSIDERED TO BE RELEVANT		
Category*	Citation of document, with indication, where appropriate, of the relevant passages	Relevant to claim No.
A	Y. GE ET AL: "flowPeaks: a fast unsupervised clustering for flow cytometry data via K-means and density peak finding", BIOINFORMATICS, vol. 28, no. 15, 1 August 2012 (2012-08-01), pages 2052-2058, XP055073112, ISSN: 1367-4803, DOI: 10.1093/bioinformatics/bts300 section 2; abstract ----- -/--	1-14
<input checked="" type="checkbox"/> Further documents are listed in the continuation of Box C. <input type="checkbox"/> See patent family annex.		
* Special categories of cited documents : "A" document defining the general state of the art which is not considered to be of particular relevance "E" earlier application or patent but published on or after the international filing date "L" document which may throw doubts on priority claim(s) or which is cited to establish the publication date of another citation or other special reason (as specified) "O" document referring to an oral disclosure, use, exhibition or other means "P" document published prior to the international filing date but later than the priority date claimed "T" later document published after the international filing date or priority date and not in conflict with the application but cited to understand the principle or theory underlying the invention "X" document of particular relevance; the claimed invention cannot be considered novel or cannot be considered to involve an inventive step when the document is taken alone "Y" document of particular relevance; the claimed invention cannot be considered to involve an inventive step when the document is combined with one or more other such documents, such combination being obvious to a person skilled in the art "&" document member of the same patent family		
Date of the actual completion of the international search 4 August 2014		Date of mailing of the international search report 11/08/2014
Name and mailing address of the ISA/ European Patent Office, P.B. 5818 Patentlaan 2 NL - 2280 HV Rijswijk Tel. (+31-70) 340-2040, Fax: (+31-70) 340-3016		Authorized officer Virnik, Elena

INTERNATIONAL SEARCH REPORT

International application No

PCT/EP2014/051963

C(Continuation). DOCUMENTS CONSIDERED TO BE RELEVANT		
Category*	Citation of document, with indication, where appropriate, of the relevant passages	Relevant to claim No.
A	DASH M ET AL: "'1+1>2': merging distance and density based clustering", DATABASE SYSTEMS FOR ADVANCED APPLICATIONS, 2001. PROCEEDINGS. SEVENTH INTERNATIONAL CONFERENCE ON APRIL 18-21, 2001, PISCATAWAY, NJ, USA, IEEE, 21 April 2001 (2001-04-21), pages 32-39, XP031977600, DOI: 10.1109/DASFAA.2001.916361 ISBN: 978-0-7695-0996-9 section 3; figure 1	1-14
A	----- K. MUMTAZ ET AL: "A Novel Density based improved k-means Clustering Algorithm - Dbkmeans", (IJCSE) INTERNATIONAL JOURNAL ON COMPUTER SCIENCE AND ENGINEERING, vol. 02, no. 02, 1 March 2010 (2010-03-01), pages 213-218, XP055074829, ISSN: 0975-3397 sections I and IV	1-14
A	----- GUNJAN GUPTA ET AL: "Bregman bubble clustering", ACM TRANSACTIONS ON KNOWLEDGE DISCOVERY FROM DATA, vol. 2, no. 2, 1 July 2008 (2008-07-01), pages 1-49, XP055074826, ISSN: 1556-4681, DOI: 10.1145/1376815.1376817 sections 1 and 3; abstract	1-14

フロントページの続き

(81) 指定国 AP(BW, GH, GM, KE, LR, LS, MW, MZ, NA, RW, SD, SL, SZ, TZ, UG, ZM, ZW), EA(AM, AZ, BY, KG, KZ, RU, TJ, T M), EP(AL, AT, BE, BG, CH, CY, CZ, DE, DK, EE, ES, FI, FR, GB, GR, HR, HU, IE, IS, IT, LT, LU, LV, MC, MK, MT, NL, NO, PL, PT, RO, R S, SE, SI, SK, SM, TR), OA(BF, BJ, CF, CG, CI, CM, GA, GN, GQ, GW, KM, ML, MR, NE, SN, TD, TG), AE, AG, AL, AM, AO, AT, AU, AZ, BA, BB, BG, BH, BN, BR, BW, BY, BZ, CA, CH, CL, CN, CO, CR, CU, CZ, DE, DK, DM, DO, DZ, EC, EE, EG, ES, FI, GB, GD, GE, GH, GM, GT, HN, H R, HU, ID, IL, IN, IR, IS, JP, KE, KG, KN, KP, KR, KZ, LA, LC, LK, LR, LS, LT, LU, LY, MA, MD, ME, MG, MK, MN, MW, MX, MY, MZ, NA, NG , NI, NO, NZ, OM, PA, PE, PG, PH, PL, PT, QA, RO, RS, RU, RW, SA, SC, SD, SE, SG, SK, SL, SM, ST, SV, SY, TH, TJ, TM, TN, TR, TT, TZ, UA, UG, US

(71) 出願人 515208636
 アンスティテュ ナシオナル ドゥ ルシエルシュ アン アンフォルマティック エ アン オ
 ートマティック
 INSTITUT NATIONAL DE RECHERCHE EN INFORMATIQ
 UE ET EN AUTOMATIQUE
 フランス共和国, エフ - 7 8 1 5 3 ル シェネ セデックス, ベーペー 1 0 5 , ロカンクール
 , ドメンヌ ドゥ ヴォリュソー

(71) 出願人 515208647
 アンスティテュ ナシオナル ドゥ ルシエルシュ アン シヤンス エ テクノロジー プール
 ランヴィロヌマン エ ラグリキュルチュール(イーエルエステーウアー)
 INSTITUT NATIONAL DE RECHERCHE EN SCIENCES E
 T TECHNOLOGIES POUR L'ENVIRONNEMENT ET L'AGR
 ICULTURE(IRSTEA)
 フランス共和国, エフ - 9 2 7 6 1 アントニ セデックス、リュ ピエール ジル ドゥ ジェ
 ンヌ セーエス 1 0 0 3 0 , 1

(74) 代理人 100080447
 弁理士 太田 恵一

(72) 発明者 セザール, ルノー
 フランス共和国, エフ - 3 0 8 7 0 クラランサック, ヴィラ 2 1 , レジデンス レ ビュイ,
 リュ デ ミコクリエ

(72) 発明者 イアンコ, ディノ
 フランス共和国, エフ - 3 4 0 9 3 モンペリエ セデックス 5 , リュ ジー . エフ プルトン
 5 0 0 , ユエムエール テーウーテーイーエス - メゾン ドゥ ラ テレデクテクシオン

(72) 発明者 マ, アンドレ
 フランス共和国, エフ - 3 4 0 9 5 モンペリエ セデックス 5 , プラス ウジェヌ パタイヨ
 ン, ユニベルシテ モンペリエ 2 , セーセー 0 5 1 , イー 3 エム

(72) 発明者 マスグリア, フロラン
 フランス共和国, エフ - 3 4 9 9 0 ジュヴィニャック, アヴニユ ドゥ カリニャン 3 1

(72) 発明者 ボンスレ, パスカル
 フランス共和国, エフ - 3 4 8 3 0 ジャク, リュ デュ ビュイ 1 8

(72) 発明者 ピュドロ, ピエール
 フランス共和国, エフ - 3 4 0 9 0 モンペリエ, リュ デ セトワヌ 7 6

(72) 発明者 スゼケリィ, エニコ
 アメリカ合衆国, ニュー ヨーク 1 0 0 0 3 , ニュー ヨーク, アパートメント 9 ケー, ワシ
 ントン プレイス 1 4

(72) 発明者 テセル, マグロンヌ
 フランス共和国, エフ - 3 4 0 9 3 モンペリエ セデックス 5 , リュ ジー . エフ プルトン
 5 0 0 , ユエムエール テーウーテーイーエス - メゾン ドゥ ラ テレデクテクシオン

(72) 発明者 ヴァンドレル, ジャン - ピエール
 フランス共和国, エフ - 3 4 1 7 0 カステルノ - ル - レ, シュマン デ マンドル 1 3 3

Fターム(参考) 2G045 BB25 CA25 CB01 CB15 FA37 FB03

专利名称(译)	如何识别罕见事件		
公开(公告)号	JP2016511397A	公开(公告)日	2016-04-14
申请号	JP2015555731	申请日	2014-01-31
[标]申请(专利权)人(译)	UNI-贝尔引用和蒙彼利埃 男性中心南某处δ统一贝尔指定Rudu蒙彼利埃 CENT HOSPITALER UNIV DE蒙彼利埃 国立研究所DUR外壳格哈德·阿南预估蝉空气污染物排放自动 法国国家信息与自动化研究所 国立研究所DUR外壳当然NSHI JANS等技术池中运行的Vie Ronu人腮格里斯的Cul薄纱庄园Ieru组阿武		
申请(专利权)人(译)	Yuniberushite蒙彼利埃 中心Osupitarie Yuniberushiteiru蒙彼利埃 研究所国立RECHERCHE安安预估蝉等安妮自动 研究所国立RECHERCHE安妮Shiyansu等技术库Ran'vironuman等滞后李的Cul薄纱(庄园Ieru组吴阿)		
[标]发明人	セザールルノー イアンコディノ マアンドレ マスグリアフロラン ポンスレパスカル ピュドロピエール スゼケリイエニコ テセルマグロンヌ ヴァンドレルジャンピエール		
发明人	セザール,ルノー イアンコ,ディノ マ,アンドレ マスグリア,フロラン ポンスレ,パスカル ピュドロ,ピエール スゼケリイ,エニコ テセル,マグロンヌ ヴァンドレル,ジャン-ピエール		
IPC分类号	G01N33/48 G01N33/53 G01N33/543 G01N33/536		
CPC分类号	G16B40/00 G01N15/1459 G01N33/569 G01N2015/1477 G01N2015/1488 G06F17/18		
FI分类号	G01N33/48.M G01N33/53.Y G01N33/53.K G01N33/543.597 G01N33/536.D		
F-TERM分类号	2G045/BB25 2G045/CA25 2G045/CB01 2G045/CB15 2G045/FA37 2G045/FB03		
代理人(译)	太田圭一		
优先权	2013153512 2013-01-31 EP		
外部链接	Espacenet		

摘要(译)

本发明提供了一种用于识别大细胞群中特定细胞的亚群，将大细胞群暴露于n种试剂并检测n种试剂的方法。通过聚类和删除非稀有细胞，将细胞聚类为k个不同的聚类。[选择图]图2

(21) 出願番号	特願2015-555731 (P2015-555731)	(71) 出願人	515121759
(86) (22) 出願日	平成26年1月31日 (2014.1.31)		ユニベルシテ ドゥ モンペリエ
(85) 翻訳文提出日	平成27年9月4日 (2015.9.4)		フランス国, 34090 モンペリエ, リ
(86) 国際出願番号	PCT/EP2014/051963		ュ オーギュスト ブルソネ 163
(87) 国際公開番号	WO2014/118343	(71) 出願人	515208625
(87) 国際公開日	平成26年8月7日 (2014.8.7)		サントル オスピタリエ ユニベルシテイ
(31) 優先権主張番号	13153512.2		ル ドゥ モンペリエ
(32) 優先日	平成25年1月31日 (2013.1.31)		CENTRE HOSPITALIER
(33) 優先権主張国	欧州特許庁 (EP)		UNIVERSITAIRE DE MO
			NTPELLIER
			フランス共和国, エフ-34295 モン
			ペリエ セデックス 5, アグニユ デュ
			ドワイヤン ガストン ジロー 191

最終頁に続く