

(19) 日本国特許庁(JP)

(12) 特 許 公 報(B2)

(11) 特許番号

特許第6160445号
(P6160445)

(45) 発行日 平成29年7月12日(2017.7.12)

(24) 登録日 平成29年6月23日(2017.6.23)

(51) Int.Cl.

F I

G 0 6 F 17/30 (2006.01)

G 0 6 F 17/30 2 1 0 D

G 0 6 F 17/30 2 2 0 Z

請求項の数 5 (全 20 頁)

(21) 出願番号 特願2013-226058 (P2013-226058)
 (22) 出願日 平成25年10月30日(2013.10.30)
 (65) 公開番号 特開2015-87966 (P2015-87966A)
 (43) 公開日 平成27年5月7日(2015.5.7)
 審査請求日 平成28年7月5日(2016.7.5)

(73) 特許権者 000005223
 富士通株式会社
 神奈川県川崎市中原区上小田中4丁目1番
 1号
 (74) 代理人 100089118
 弁理士 酒井 宏明
 (72) 発明者 松本 和宏
 神奈川県川崎市中原区上小田中4丁目1番
 1号 富士通株式会社内
 審査官 齊藤 貴孝

最終頁に続く

(54) 【発明の名称】 分析装置、分析方法および分析プログラム

(57) 【特許請求の範囲】

【請求項1】

入力データに対してサンプリングを実行し前記入力データから一部のデータを抽出する処理を繰り返し実行して複数のサンプリングデータを生成するサンプリング実行部と、

前記複数のサンプリングデータについてクラスタ分析を実行し、前記サンプリングデータ毎に、前記サンプリングデータに含まれるデータを異なるクラスタに分類するクラスタ分析部と、

前記複数のサンプリングデータに対する前記クラスタ分析部の複数の分類結果と前記入力データとを基にして、前記入力データに含まれるデータの所属するクラスタを予測したデータを示す予測データを複数生成するクラスタ予測部と、

前記予測データのクラスタ間距離およびクラスタ内距離を基にして、予測データ毎に評価値を算出し、パレート解となる評価値に対応する予測データを判定する判定部と、

前記パレート解となる評価値に対応する予測データを基にして、前記入力データに含まれるデータをクラスタに分類する最終クラスタ計算部と

を有することを特徴とする分析装置。

【請求項2】

前記最終クラスタ計算部は、パレート解となる評価値に対応する類似の予測データをグループ化し、同一グループに含まれる予測データを基にして、前記入力データに含まれるデータを異なるクラスタに分類する処理を、グループ毎に実行することを特徴とする請求項1に記載の分析装置。

【請求項 3】

前記最終クラスタ計算部は、前記入力データに対して、ランダムにクラスタを割り当てた複数の最終クラスタデータを生成し、各最終クラスタデータと予測データとの類似度を基にして、特定の最終クラスタデータを選択することを特徴とする請求項 1 または 2 に記載の分析装置。

【請求項 4】

コンピュータが実行する分析方法であって、

入力データに対してサンプリングを実行し前記入力データから一部のデータを抽出する処理を繰り返し実行して複数のサンプリングデータを生成し、

前記複数のサンプリングデータについてクラスタ分析を実行し、前記サンプリングデータ毎に、前記サンプリングデータに含まれるデータを異なるクラスタに分類し、

前記複数のサンプリングデータに対する前記クラスタ分析部の複数の分類結果と前記入力データとを基にして、前記入力データに含まれるデータの所属するクラスタを予測したデータを示す予測データを複数生成し、

前記予測データのクラスタ間距離およびクラスタ内距離を基にして、予測データ毎に評価値を算出し、パレート解となる評価値に対応する予測データを判定し、

前記パレート解となる評価値に対応する予測データを基にして、前記入力データに含まれるデータをクラスタに分類する

各処理を実行することを特徴とする分析方法。

【請求項 5】

コンピュータに、

入力データに対してサンプリングを実行し前記入力データから一部のデータを抽出する処理を繰り返し実行して複数のサンプリングデータを生成し、

前記複数のサンプリングデータについてクラスタ分析を実行し、前記サンプリングデータ毎に、前記サンプリングデータに含まれるデータを異なるクラスタに分類し、

前記複数のサンプリングデータに対する前記クラスタ分析部の複数の分類結果と前記入力データとを基にして、前記入力データに含まれるデータの所属するクラスタを予測したデータを示す予測データを複数生成し、

前記予測データのクラスタ間距離およびクラスタ内距離を基にして、予測データ毎に評価値を算出し、パレート解となる評価値に対応する予測データを判定し、

前記パレート解となる評価値に対応する予測データを基にして、前記入力データに含まれるデータをクラスタに分類する

各処理を実行させることを特徴とする分析プログラム。

【発明の詳細な説明】**【技術分野】****【0001】**

本発明は、分析装置等に関する。

【背景技術】**【0002】**

クラスタ分析は、データの集まりをデータ間の類似度に基づいて複数のクラスタに分類する処理である。例えば、クラスタ分析には、階層的クラスタ分析や非階層的クラスタ分析がある。

【0003】

階層的クラスタ分析は、例えば、個々のデータを 1 つのクラスタとして設定し、クラスタ間の類似度を計算し、最も類似している各クラスタを併合する処理を繰り返し実行するものである。

【0004】

非階層的クラスタ分析は、分類の状態を表す関数を使い、関数の値が最適解となるように探索を行うものである。

【先行技術文献】

【特許文献】

【0005】

【特許文献1】特開2007-179143号公報

【特許文献2】特開2005-293048号公報

【発明の概要】

【発明が解決しようとする課題】

【0006】

しかしながら、上述した従来技術では、大規模データについてクラスタ分析を実行すると時間を要するという問題がある。

【0007】

例えば、階層的クラスタ分析および非階層的クラスタ分析はそれぞれ、小規模データ、中規模データに対してクラスタ分析を実行することを想定している。このため、現実的な計算機環境により、大規模データに対して階層的クラスタ分析や非階層的クラスタ分析を実行すると、現実的な計算時間内で計算できないことがある。

【0008】

1つの側面では、クラスタ分析に要する時間を削減することができる分析装置、分析方法および分析プログラムを提供することを目的とする。

【課題を解決するための手段】

【0009】

第1の案では、分析装置は、サンプリング実行部、クラスタ分析部、クラスタ予測部、判定部、最終クラスタ計算部を有する。サンプリング実行部は、入力データに対してサンプリングを実行し前記入力データから一部のデータを抽出する処理を繰り返し実行して複数のサンプリングデータを生成する。クラスタ分析部は、複数のサンプリングデータについてクラスタ分析を実行し、サンプリングデータ毎に、サンプリングデータに含まれるデータを異なるクラスタに分類する。クラスタ予測部は、複数のサンプリングデータに対するクラスタ分析部の複数の分類結果と入力データとを基にして、入力データに含まれるデータの所属するクラスタを予測したデータを示す予測データを複数生成する。判定部は、予測データのクラスタ間距離およびクラスタ内距離を基にして、予測データ毎に評価値を算出し、パレート解となる評価値に対応する予測データを判定する。最終クラスタ計算部は、パレート解となる評価値に対応する予測データを基にして、前記入力データに含まれるデータをクラスタに分類する。なお、パレート解となる予測データは、例えば、評価値が他の予測データと比較して優越するものである。

【発明の効果】

【0010】

本発明の1実施態様によれば、クラスタ分析に要する時間を削減することができるという効果を奏する。

【図面の簡単な説明】

【0011】

【図1】図1は、本実施例1に係る分析装置の構成を示す機能ブロック図である。

【図2】図2は、分析対象データのデータ構造の一例を示す図である。

【図3】図3は、サンプリングデータテーブルのデータ構造の一例を示す図である。

【図4】図4は、予測データテーブルのデータ構造の一例を示す図である。

【図5】図5は、評価値データテーブルのデータ構造の一例を示す図である。

【図6】図6は、中間データテーブルのデータ構造の一例を示す図である。

【図7】図7は、最終データのデータ構造の一例を示す図である。

【図8】図8は、各予測データのクラスタ内距離とクラスタ間距離との関係を示す図(1)である。

【図9】図9は、最終データ候補テーブルの一例を示す図である。

【図10】図10は、本実施例1にかかる分析装置の処理手順を示すフローチャートである。

10

20

30

40

50

【図 1 1】図 1 1 は、本実施例 2 にかかる分析装置の構成を示す機能ブロック図である。

【図 1 2】図 1 2 は、各予測データのクラスタ内距離とクラスタ間距離との関係を示す図 (2) である。

【図 1 3】図 1 3 は、実施例 2 にかかる分析装置の処理手順を示すフローチャート (1) である。

【図 1 4】図 1 4 は、実施例 2 にかかる分析装置の処理手順を示すフローチャート (2) である。

【図 1 5】図 1 5 は、分析プログラムを実行するコンピュータの一例を示す図である。

【発明を実施するための形態】

【 0 0 1 2 】

10

以下に、本願の開示する分析装置、分析方法および分析プログラムの実施例を図面に基
づいて詳細に説明する。なお、この実施例によりこの発明が限定されるものではない。

【実施例 1】

【 0 0 1 3 】

図 1 は、本実施例 1 に係る分析装置の構成を示す機能ブロック図である。図 1 に示す
ように、この分析装置 1 0 0 は、通信部 1 1 0、入力部 1 2 0、出力部 1 3 0、記憶部 1 4
0、制御部 1 5 0 を有する。

【 0 0 1 4 】

通信部 1 1 0 は、無線または有線によってネットワークに接続し、ネットワークを介し
て、他の装置とデータ通信を行う処理部である。通信部 1 1 0 は、通信装置に対応する。

20

【 0 0 1 5 】

入力部 1 2 0 は、各種の情報を入力する入力装置である。入力部 1 2 0 は、例えば、キ
ーボードやマウス、タッチパネル等に対応する。

【 0 0 1 6 】

出力部 1 3 0 は、制御部 1 5 0 から出力される情報を表示する表示装置である。例えば
、出力部 1 3 0 は、モニタ、液晶ディスプレイ、タッチパネル等に対応する。

【 0 0 1 7 】

記憶部 1 4 0 は、分析対象データ 1 4 1、サンプリングデータテーブル 1 4 2、予測デ
ータテーブル 1 4 3、評価値テーブル 1 4 4、中間データテーブル 1 4 5、最終データ 1
4 6 を有する。記憶部 1 4 0 は、例えば、R A M (Random Access Memory)、R O M (30
Read Only Memory)、フラッシュメモリ (Flash Memory) などの半導体メモリ素子や
、H D D (Hard Disk Drive) などの記憶装置に対応する。

30

【 0 0 1 8 】

分析対象データ 1 4 1 は、クラスタ分析の対象となるデータである。図 2 は、分析対象
データのデータ構造の一例を示す図である。図 2 に示すように、分析対象データ 1 4 1 は
、識別番号、年齢、性別、身長、体重等を有する。識別番号は、各レコードを一意に識別
する情報である。年齢、性別、身長、体重は、特定の人物の年齢、性別、身長、体重をそ
れぞれ示す情報である。なお、図 2 に示す例では、性別を 1 または 2 で表す。例えば、性
別「 1 」は、性別が男性であることを示し、性別「 2 」は、性別が女性であることを示
す。

40

【 0 0 1 9 】

サンプリングデータテーブル 1 4 2 は、複数のサンプリングデータを有するテーブルで
ある。各サンプリングデータは、後述するサンプリング実行部 1 5 1 によって生成される
。サンプリング実行部 1 5 1 が分析対象データ 1 4 1 をサンプリングすることで、各サン
プリングデータが生成される。図 3 は、サンプリングデータテーブルのデータ構造の一例
を示す図である。図 3 に示すように、サンプリングデータテーブル 1 4 2 は、サンプリ
ングデータ 1 4 2 a , 1 4 2 b , 1 4 2 c を有する。図 3 では一例として、サンプリングデ
ータ 1 4 2 a , 1 4 2 b , 1 4 2 c を示すが、その他のサンプリングデータを含んでも良
い。

【 0 0 2 0 】

50

図3において、例えば、サンプリングデータ142aは、識別番号、年齢、性別、身長、体重、クラスタ番号を有する。識別番号、年齢、性別、身長、体重に関する説明は、図2で説明した、年齢、性別、身長、体重の説明と同様である。

【0021】

予測データテーブル143は、複数の予測データを有するテーブルである。各予測データは、後述するクラスタ予測部153によって生成される。クラスタ予測部153が、サンプリングデータを基にして、分析対象データ141の各レコードのクラスタ番号を予測することで、予測データを生成する。サンプリングデータ毎に予測データが生成される。図4は、予測データテーブルのデータ構造の一例を示す図である。図4に示すように、予測データテーブル143は、予測データ143a, 143b, 143cを有する。図4では一例として、予測データ143a, 143b, 143cを示すが、その他の予測データを含んでも良い。

10

【0022】

評価値データテーブル144は、各予測データの評価値をそれぞれ保持するテーブルである。図5は、評価値データテーブルのデータ構造の一例を示す図である。図5に示すように、この評価値データテーブル144は、予測データ識別情報と、評価値とを対応付ける。予測データ識別情報は、予測データを一意に識別する情報である。

【0023】

図5において、評価値は、クラスタ間距離と、クラスタ内距離とを含む。クラスタ間距離は、異なるクラスタ間の距離を示すものである。一般的に、クラスタ間距離が大きいほど、クラスタ分析結果に対する評価が高くなる。クラスタ内距離は、クラスタの直径を示すものである。一般的に、クラスタ内距離が小さいほど、クラスタ分析結果に対する評価が高くなる。即ち、クラスタ間距離が大きいほど、また、クラスタ内距離が小さいほど、クラスタ分析結果が優れている。

20

【0024】

中間データテーブル145は、複数の中間データを有するテーブルである。各中間データは、評価値の良い予測データに対応して作成される。図6は、中間データテーブルのデータ構造の一例を示す図である。図6に示すように、この中間データテーブル145は、中間データ145a, 145g, 145zを有する。図6では一例として、中間データ145a, 145g, 145zを示すが、その他の中間データを含んでも良い。

30

【0025】

図6において、中間データは、識別番号と、各クラスタ番号とを対応付ける。識別番号は、分析対象データ141の識別番号に対応する。例えば、中間データ145aの1段目では、識別番号1001に対応するレコードが、クラスタ番号「1」に分類されることを示す。ここで、識別番号「1001」に対応するレコードは、図2に示した分析対象データ141の識別番号「1001」に対応するレコードに対応する。従って、図2に示した分析対象データ141の識別番号「1001」に対応するレコードが、クラスタ番号「1」のクラスタに属していることを示す。

【0026】

最終データ146は、分析対象データ141の最終的なクラスタ分析結果を示す。図7は、最終データのデータ構造の一例を示す図である。図7に示すように、この最終データ146は、識別番号と、各クラスタ番号とを対応付ける。識別番号は、分析対象データ141の識別番号に対応する。例えば、最終データ146の1段目では、識別番号1001に対応するレコードが、クラスタ番号「1」に分類されることを示す。

40

【0027】

図1の説明に戻る。制御部150は、サンプリング実行部151、クラスタ分析部152、クラスタ予測部153、判定部154、最終クラスタ計算部155を有する。制御部150は、例えば、ASIC (Application Specific Integrated Circuit) や、FPGA (Field Programmable Gate Array) などの集積装置に対応する。また、制御部150は、例えば、CPU (Central Processing Unit) やMPU (Micro Processing

50

Unit)等の電子回路に対応する。

【0028】

サンプリング実行部151は、分析対象データ141に対してサンプリングを複数回繰り返し実行することで、複数のサンプリングデータを生成する処理部である。サンプリング実行部151は、生成した各サンプリングデータを、サンプリングデータテーブル142に格納する。

【0029】

例えば、サンプリング実行部151は、入力部120を介して、計算回数およびサンプリング件数を取得し、取得した計算回数だけ、サンプリングを行う。また、サンプリング実行部151は、サンプリングを実行する度に、サンプリング間隔を変更しても良い。また、ランダムサンプリングを行っても良い。

10

【0030】

サンプリング実行部151は、サンプリングデータのレコードの件数を、入力部120から取得したサンプリング件数に合わせる。例えば、指定されたサンプリング件数がN2件の場合には、各サンプリングデータの件数をそれぞれN2件とする。例えば、分析対象データのレコードの件数をN1件とすると、N1とN2との大小関係は「 $N1 > N2$ 」となる。

【0031】

クラスタ分析部152は、サンプリングデータテーブル142に格納された各サンプリングデータを取得し、各サンプリングデータをクラスタ分析する処理部である。クラスタ分析部152は、クラスタ分析結果に応じて、サンプリングデータの各レコードについてクラスタ番号を割り当てる。

20

【0032】

図3に示したサンプリングデータテーブル142を例にして説明を行う。クラスタ分析部152は、まず、サンプリングデータ142aに対してクラスタ分析を行い、サンプリングデータ142aの各レコードを複数のクラスタに分類し、分類結果に応じて、クラスタ番号を割り振る。クラスタ分析部152は、サンプリングデータ142b, 142cについても同様に、クラスタ分析を行って、各レコードを、複数のクラスタに分類し、分類結果に応じて、クラスタ番号を割り振る。クラスタ分析部152が分類するクラスタの数は、予め設定されているものとする。

30

【0033】

クラスタ分析部152が行うクラスタ分析は、階層的クラスタ分析でも良いし、非階層的クラスタ分析でもよい。ここでは一例として、クラスタ分析部152が、階層的クラスタ分析を実行する場合について説明する。

【0034】

クラスタ分析部152が、階層的クラスタ分析を行う場合には、まず、個々のデータを1つのクラスタとして設定し、クラスタ間の類似度を計算する。クラスタ分析部152は、最も類似しているクラスタを併合する。クラスタ分析部152は、予め設定されたクラスタの数と同数になるまで、上記処理を繰り返し実行する。

【0035】

例えば、クラスタ分析部152は、各クラスタの組み合わせについて、クラスタ間のユークリッド距離を算出し、ユークリッド距離が最小となる各クラスタの組みを、合併する。この場合には、クラスタ間のユークリッド距離が上記クラスタ間の類似度に対応し、ユークリッド距離が短いほど、類似度が高い。

40

【0036】

クラスタ予測部153は、サンプリングデータテーブル142のサンプリングデータのクラスタ分析結果に基づいて、分析対象データ141の各レコードのクラスタ番号を予測し、予測データテーブル143を生成する処理部である。クラスタ予測部153は、サンプリングデータテーブル142に含まれるサンプリングデータの数だけ、予測データを生成し、生成した予測データを予測データテーブル143に登録する。

50

【 0 0 3 7 】

例えば、クラスタ予測部 1 5 3 は、サンプリングデータテーブル 1 4 2 のサンプリングデータ 1 4 2 a を基にして、予測データ 1 4 3 a を生成する。クラスタ予測部 1 5 3 は、サンプリングデータ 1 4 2 b を基にして、予測データ 1 4 3 b を生成する。クラスタ予測部 1 5 3 は、サンプリングデータ 1 4 2 c を基にして、予測データ 1 4 3 c を生成する。クラスタ予測部 1 5 3 は、サンプリングデータが N 個存在する場合には、予測データを N 個作成する。

【 0 0 3 8 】

ここで、クラスタ予測部 1 5 3 が、サンプリングデータ 1 4 2 a を基にして、予測データ 1 4 3 a を生成する場合の処理の一例について説明する。まず、クラスタ予測部 1 5 3 は、サンプリングデータ 1 4 2 に含まれる識別番号と、クラスタ番号との関係を、そのまま、予測データ 1 4 3 a に設定する。

10

【 0 0 3 9 】

例えば、クラスタ予測部 1 5 3 は、サンプリングデータ 1 4 2 a に識別番号「1 0 0 1」のレコードのクラスタ番号が「1」の場合には、予測データ 1 4 3 a の識別番号「1 0 0 1」のクラスタ番号を「1」に設定する。同様に、クラスタ予測部 1 5 3 は、サンプリングデータ 1 4 2 a に存在する全ての識別番号とクラスタ番号との関係を、予測データ 1 4 3 a に設定する。

【 0 0 4 0 】

続いて、クラスタ予測部 1 5 3 は、上記処理を行った結果、クラスタ番号が未設定となるレコードについて下記の処理を行う。まず、クラスタ予測部 1 5 3 は、各クラスタに分類されたレコードから、代表レコードを検出する。例えば、クラスタ番号「1」のレコードのうち、平均的な数値を有するレコードを代表レコードとして検出する。クラスタ予測部 1 5 3 は、他のクラスタ番号に対応する代表レコードも同様に検出する。

20

【 0 0 4 1 】

クラスタ予測部 1 5 3 は、クラスタ番号が未設定のレコードと、各代表レコードとのユークリッド距離を計算し、ユークリッド距離が最小となる組み合わせを特定する。クラスタ予測部 1 5 3 は、特定した組の代表レコードのクラスタ番号を、該当するレコードのクラスタ番号に設定する。

【 0 0 4 2 】

例えば、クラスタ番号が未設定のレコードと、各代表レコードとのユークリッド距離を算出し、未設定のレコードと、クラスタ番号「1」の代表レコードとのユークリッド距離が最小の場合には、該当するレコードのクラスタ番号を「1」に設定する。クラスタ予測部 1 5 3 は、未設定のレコードについて、上記処理を繰り返し実行することで、予測データテーブル 1 4 3 を生成する。

30

【 0 0 4 3 】

判定部 1 5 4 は、予測データテーブル 1 4 3 を基にして、評価値データテーブル 1 4 4 を生成する処理部である。評価部 1 5 4 は、予測データテーブル 1 4 3 に含まれる予測データ毎に評価値を算出する。

【 0 0 4 4 】

判定部 1 5 4 は、予測データ毎にクラスタ間距離およびクラスタ内距離を算出し、クラスタ間距離およびクラスタ内距離を予測データの評価値とする。予測データのクラスタ間距離を算出する処理の一例について説明する。ここでは、クラスタ番号「1 ~ 3」のクラスタが存在するものとする。判定部 1 5 4 は、クラスタ番号「1」に属する第 1 代表レコードと、クラスタ番号「2」に属する第 2 代表レコードと、クラスタ番号「3」に属する第 3 代表レコードとを検出する。代表レコードを検出する処理の一例は、例えば、同一のクラスタ番号に属するレコードのうち、平均的な数値を有するレコードを代表レコードとして検出する。

40

【 0 0 4 5 】

判定部 1 5 4 は、第 1 代表レコードと、第 2 代表レコードとのユークリッド距離を算出

50

し、第1代表レコードと第3代表レコードとのユークリッド距離を算出する。判定部154は、算出した各ユークリッド距離を平均したユークリッド距離を、予測データのクラスタ間距離とする。

【0046】

例えば、第1代表レコードの年齢、性別、身長、体重の値をそれぞれ、a1、a2、a3、a4とする。第2代表レコードの年齢、性別、身長、体重の値をそれぞれ、b1、b2、b3、b4とする。第3代表レコードの年齢、性別、身長、体重の値をそれぞれ、c1、c2、c3、c4とする。この場合には、第1代表レコードと、第2代表レコードとのユークリッド距離X1は、式(1)で計算され、第1代表レコードと、第3代表レコードとのユークリッド距離X2は、式(2)で計算される。この場合には、予測データのクラスタ間距離は式(3)に示すものとなる。

10

【0047】

ユークリッド距離 $X1 = ((a1 - b1)^2 + (a2 - b2)^2 + (a3 - b3)^2 + (a4 - b4)^2)^{1/2} \dots (1)$

【0048】

ユークリッド距離 $X2 = ((a1 - c1)^2 + (a2 - c2)^2 + (a3 - c3)^2 + (a4 - c4)^2)^{1/2} \dots (2)$

【0049】

クラスタ間距離 $= (X1 + X2) / 2 \dots (3)$

【0050】

20

続いて、クラスタ内距離を算出する処理について説明する。まず、判定部154は、同一のクラスタ番号に属する各レコード間のユークリッド距離をそれぞれ算出する。そして、判定部154は、算出したユークリッド距離を平均したユークリッド距離を、予測データのクラスタ内距離とする。判定部154は、各クラスタ番号のクラスタに対応するクラスタ内距離を平均することで、予測データのクラスタ内距離を算出する。

【0051】

例えば、クラスタ番号「1～3」のクラスタが存在する場合には、判定部154は、各クラスタ番号「1～3」のクラスタ内距離をそれぞれ算出する。判定部154は、各クラスタ番号「1～3」のクラスタ内距離を平均することで、予測データのクラスタ内距離を算出する。

30

【0052】

例えば、クラスタ番号「1」のクラスタ内距離を算出する例について説明する。クラスタ内に3つの第1レコード、第2レコード、第3レコードが存在するものとする。例えば、第1レコードの年齢、性別、身長、体重の値をそれぞれ、d1、d2、d3、d4とする。第2レコードの年齢、性別、身長、体重の値をそれぞれ、e1、e2、e3、e4とする。第3代表レコードの年齢、性別、身長、体重の値をそれぞれ、f1、f2、f3、f4とする。この場合には、第1レコードと、第2レコードとのユークリッド距離Y1は、式(4)で計算され、第1レコードと、第3レコードとのユークリッド距離Y2は、式(5)で計算される。この場合には、クラスタ番号「1」のクラスタのクラスタ内距離は式(6)に示すものとなる。

40

【0053】

ユークリッド距離 $Y1 = ((d1 - e1)^2 + (d2 - e2)^2 + (d3 - e3)^2 + (d4 - e4)^2)^{1/2} \dots (4)$

【0054】

ユークリッド距離 $Y2 = ((d1 - f1)^2 + (d2 - f2)^2 + (d3 - f3)^2 + (d4 - f4)^2)^{1/2} \dots (5)$

【0055】

クラスタ内距離 $= (Y1 + Y2) / 2 \dots (6)$

【0056】

判定部154は、他のクラスタについても同様にクラスタ内距離を算出し、各クラスタ

50

のクラスタ内距離を平均することで、予測データのクラスタ内距離を算出する。

【 0 0 5 7 】

判定部 1 5 4 は、予測データテーブル 1 4 3 に含まれる予測データ毎に上記処理を実行することで、各予測データの評価値を算出し、評価値データテーブル 1 4 4 を生成する。

【 0 0 5 8 】

最終クラスタ計算部 1 5 5 は、分析対象データ 1 4 1 の最終的なクラスタ分析結果となる最終データ 1 4 6 を生成する処理部である。最終クラスタ計算部 1 5 5 は、評価値データテーブル 1 4 4 から中間データテーブル 1 4 5 を生成する処理を行った後に、中間データテーブル 1 4 5 を基にして、最終データ 1 4 6 を生成する。

【 0 0 5 9 】

最終クラスタ計算部 1 5 5 が、評価値データテーブル 1 4 4 から中間データテーブル 1 4 5 を生成する処理の一例について説明する。最終クラスタ計算部 1 5 5 は、評価値データテーブル 1 4 4 の予測データ毎の評価値を比較して、パレート解となる予測データを特定し、特定したパレート解となる予測データを、中間データテーブル 1 4 5 に設定する。例えば、パレート解となる予測データは、一つ以上の項目について他の予測データよりも優れているものとなる。

【 0 0 6 0 】

図 8 は、各予測データのクラスタ内距離とクラスタ間距離との関係を示す図 (1) である。図 8 において、縦軸はクラスタ内距離を示し、横軸はクラスタ間距離を示す。一般的に、クラスタ間距離が大きいほど、また、クラスタ内距離が小さいほど、予測データは、
20

【 0 0 6 1 】

続いて、最終クラスタ計算部 1 5 5 が、中間データテーブル 1 4 5 から最終データ 1 4 6 を生成する処理について説明する。まず、最終クラスタ計算部 1 5 5 は、最終データ候補テーブルを生成する。図 9 は、最終データ候補テーブルの一例を示す図である。図 9 に示すように、この最終データ候補テーブル 1 0 は、最終データ候補 1 0 a , 1 0 b , 1 0 c を有する。ここでは一例として、最終データ候補 1 0 a , 1 0 b , 1 0 c を示すが、これ以外に、最終データ候補を含んでいても良い。

【 0 0 6 2 】

最終クラスタ計算部 1 5 5 は、最終データ候補 1 0 a , 1 0 b , 1 0 c の各クラスタ番号を 0 の初期値に設定する。そして、最終クラスタ計算部 1 5 5 は、各識別番号の各クラスタ番号の値のいずれか一つが「 1 」となるように、ランダムに「 1 」を割り振る。例えば、図 9 に示す例では、最終データ候補 1 0 a の識別番号「 1 0 0 1 」に対してランダムに「 1 」を割り振ることで、クラスタ番号「 1 」に対応するものが「 1 」に設定され、その他のクラスタ番号については「 0 」が設定される。

【 0 0 6 3 】

最終クラスタ計算部 1 5 5 は、最終データ候補テーブル 1 0 の各最終データ候補 1 0 a , 1 0 b , 1 0 c と、中間データテーブル 1 4 5 の各中間データとの類似度を計算し、最も類似度の高い最終データ候補を、最終データ 1 4 6 として特定する。
40

【 0 0 6 4 】

最終クラスタ計算部 1 5 5 は、中間データの識別番号および識別番号に対応するクラスタ番号と、最終データ候補の識別番号および識別番号に対応するクラスタ番号とを比較し、一致する数を計数する。最終クラスタ計算部 1 5 5 は、一致する数を、全レコード数で除算することで、類似度を算出する。以下の説明では、一致する数を、一致数と表記する。

【 0 0 6 5 】

例えば、最終クラスタ計算部 1 5 5 が、最終データ候補 1 0 a の類似度を算出する場合について説明する。最終クラスタ計算部 1 5 5 は、最終データ候補 1 0 a と中間データ 1 4 5 a とを比較し、一致数が「 L 1 」であり、最終データ候補 1 0 a の全レコード数が「
50

M1」の場合には、最終データ候補10aと中間データ145aとの類似度は「 $L1/M1$ 」となる。最終クラスタ計算部155は、最終データ候補10aと中間データ145gとを比較し、一致数が「L2」であり、最終データ候補10aの全レコード数が「M2」の場合には、最終データ候補10aと中間データ145aとの類似度は「 $L2/M2$ 」となる。最終クラスタ計算部155は、最終データ候補10aと中間データ145zとを比較し、一致数が「L3」であり、最終データ候補10aの全レコード数が「M3」の場合には、最終データ候補10aと中間データ145zとの類似度は「 $L3/M3$ 」となる。この場合には、最終クラスタ計算部155は、最終データ候補10aの類似度を「 $L1/M1 + L2/M2 + L3/M3$ 」と特定する。

【0066】

最終クラスタ計算部155は、最終データ候補10b, 10cに関しても、最終データ候補10aと同様にして、類似度を算出する。最終クラスタ計算部155は、最終データ候補10aの類似度、最終データ候補10bの類似度、最終データ候補10cの類似度を比較し、類似度が最大となる最終データ候補を特定する。最終クラスタ計算部155は、特定した最終データ候補を、最終データ146として設定する。最終クラスタ計算部155は、最終データ146を、出力部130に出力しても良い。

【0067】

次に、本実施例1にかかる分析装置100の処理手順について説明する。図10は、本実施例1にかかる分析装置の処理手順を示すフローチャートである。図10に示すように、分析装置100は、分析対象データ141を受け付ける（ステップS101）。また、分析装置100は、繰り返し計算回数を受け付け（ステップS102）、サンプリング件数を受け付ける。また、カウント値を初期化する（ステップS103）。分析装置100は、カウント値に1を加算する（ステップS104）。カウント値の初期値を0とする。

【0068】

分析装置100は、分析対象データ141をサンプリングし、サンプリングデータを生成する（ステップS105）。各サンプリングデータは、サンプリングデータテーブル142に格納される。分析装置100は、サンプリングデータに対してクラスタ分析処理を実行し、各々のレコードに対してクラスタ番号を割り振る（ステップS106）。

【0069】

分析装置100は、クラスタ番号を割り振ったサンプリングデータと分析対象データ141とを比較して、分析対象データ141に含まれる各々のレコードに対してクラスタ番号を割り振ることで予測データを生成する（ステップS107）。各予測データは、予測データテーブル143に格納される。

【0070】

分析装置100は、予測データを基にして、クラスタ内距離およびクラスタ間距離を算出し、評価値データテーブル144を生成する（ステップS108）。分析装置100は、繰り返しの計算回数がカウント値未満であるか否かを判定する（ステップS109）。分析装置100は、繰り返しの計算回数がカウント値未満の場合には（ステップS1090, Yes）、ステップS104に移行する。

【0071】

一方、分析装置100は、繰り返しの計算回数がカウント値以上である場合には（ステップS109, No）、パレート解に対応する予測データを選択して、中間データテーブル145を作成する（ステップS110）。

【0072】

分析装置100は、ランダムにクラスタ番号を割り振った複数の最終データ候補を生成する（ステップS111）。ステップS111において、分析装置100は、類似度が大きくなるようにクラスタ番号を割り振る。例えば、分析装置100は、ランダムにクラスタ番号を割り振り、類似度を計算する。そして、分析装置100は、類似度が大きい、クラスタ番号の割り振りを少し変更して、類似度が大きくなるか、試行する処理を利用者が設定した回数繰り返す。

10

20

30

40

50

【 0 0 7 3 】

分析装置 1 0 0 は、中間データと各最終データ候補とを比較して、類似度が最大となる最終データ候補を判定する（ステップ S 1 1 2）。ステップ S 1 1 2 で判定した類似度が最大となる最終データ候補が、最終データ 1 4 6 となる。分析装置 1 0 0 は、判定結果を出力する（ステップ S 1 1 3）。

【 0 0 7 4 】

次に、本実施例 1 にかかる分析装置 1 0 0 の効果について説明する。分析装置 1 0 0 は分析対象データ 1 4 1 から抽出したサンプリングデータをクラスタ分析し、サンプリングデータのクラスタ分析結果を基にして、分析対象データの各データが属するクラスタを予測した複数の予測データを生成する。そして、分析装置 1 0 0 は、複数の予測データのうち、評価値のよい予測データのクラスタ分類結果を用いて、分析対象データ 1 4 1 の最終的なクラスタ分類結果を特定する。これにより、分析装置 1 0 0 によれば、クラスタ分析に要する時間を削減することができる。

10

【 0 0 7 5 】

また、現実的な計算機で、現実的な時間内に計算できない、大規模なデータに対して、現実的な計算機で、現実的な時間内に、クラスタ分析を実行することができる。

【実施例 2】

【 0 0 7 6 】

図 1 1 は、本実施例 2 にかかる分析装置の構成を示す機能ブロック図である。図 1 1 に示すように、この分析装置 2 0 0 は、通信部 2 1 0、入力部 2 2 0、出力部 2 3 0、記憶部 2 4 0、制御部 2 5 0 を有する。

20

【 0 0 7 7 】

通信部 2 1 0、入力部 2 2 0、出力部 2 3 0 に関する説明は、図 1 に示した、通信部 1 1 0、入力部 1 2 0、出力部 1 3 0 に関する説明と同様である。

【 0 0 7 8 】

記憶部 2 4 0 は、分析対象データ 2 4 1、サンプリングデータテーブル 2 4 2、予測データテーブル 2 4 3、評価値データテーブル 2 4 4、中間データテーブル 2 4 5、最終データ 2 4 6 を有する。記憶部 2 4 0 は、例えば、RAM、ROM、フラッシュメモリなどの半導体メモリ素子や、HDD などの記憶装置に対応する。

30

【 0 0 7 9 】

分析対象データ 2 4 1 は、クラスタ分析の対象となるデータである。分析対象データ 2 4 1 のデータ構造は、図 2 に示した分析対象データ 1 4 1 のデータ構造と同様である。

【 0 0 8 0 】

サンプリングデータテーブル 2 4 2 は、複数のサンプリングデータを有するテーブルである。各サンプリングデータは、後述するサンプリング実行部 2 5 1 によって生成される。サンプリング実行部 2 5 1 が分析対象データ 2 4 1 をサンプリングすることで、各サンプリングデータが生成される。サンプリングデータテーブル 2 4 2 のデータ構造は、図 3 に示したサンプリングデータテーブル 1 4 2 のデータ構造と同様である。

【 0 0 8 1 】

予測データテーブル 2 4 3 は、複数の予測データを有するテーブルである。各予測データは、後述するクラスタ予測部 2 5 3 によって生成される。クラスタ予測部 2 5 3 が、サンプリングデータを基にして、分析対象データ 2 4 1 の各レコードのクラスタ番号を予測することで、予測データを生成する。サンプリングデータ毎に予測データが生成される。予測データテーブル 2 4 3 のデータ構造は、図 4 に示した予測データテーブル 1 4 3 のデータ構造と同様である。

40

【 0 0 8 2 】

評価値データテーブル 2 4 4 は、各予測データの評価値をそれぞれ保持するテーブルである。評価値データテーブル 2 4 4 のデータ構造は、図 5 に示した評価値データテーブル 1 4 4 のデータ構造と同様である。

【 0 0 8 3 】

50

中間データテーブル 245 は、複数の中間データを有するテーブルである。各中間データは、評価値の良い予測データに対応して作成される。中間データテーブル 245 のデータ構造は、図 6 に示した中間データテーブル 145 のデータ構造と同様である。

【0084】

最終データ 246 a , 246 b , 246 c は、分析対象データ 241 の最終的なクラスタ分析結果を示す。各最終データ 246 a , 246 b , 246 c のデータ構造は、図 7 に示した最終データ 146 のデータ構造と同様である。

【0085】

図 11 の説明に戻る。制御部 250 は、サンプリング実行部 251、クラスタ分析部 252、クラスタ予測部 253、判定部 254、最終クラスタ計算部 255 を有する。制御部 250 は、例えば、ASIC や、FPGA などの集積装置に対応する。また、制御部 250 は、例えば、CPU や MPU 等の電子回路に対応する。

10

【0086】

サンプリング実行部 251 は、分析対象データ 241 に対してサンプリングを複数回繰り返し実行することで、複数のサンプリングデータを生成する処理部である。サンプリング実行部 251 は、生成した各サンプリングデータを、サンプリングデータテーブル 242 に格納する。サンプリング実行部 251 の具体的な処理は、図 1 に示したサンプリング実行部 151 と同様である。

【0087】

クラスタ分析部 252 は、サンプリングデータテーブル 242 に格納された各サンプリングデータを取得し、各サンプリングデータをクラスタ分析する処理部である。クラスタ分析部 252 は、クラスタ分析結果に応じて、サンプリングデータの各レコードについてクラスタ番号を割り当てる。クラスタ分析部 252 の具体的な処理は、図 1 に示したクラスタ分析部 152 と同様である。

20

【0088】

クラスタ予測部 253 は、サンプリングデータテーブル 242 のサンプリングデータのクラスタ分析結果に基づいて、分析対象データ 241 の各レコードのクラスタ番号を予測し、予測データテーブル 243 を生成する処理部である。クラスタ予測部 253 の具体的な処理は、図 1 に示したクラスタ予測部 153 と同様である。

【0089】

判定部 254 は、予測データテーブル 243 を基にして、評価値データテーブル 244 を生成する処理部である。判定部 254 は、予測データテーブル 243 に含まれる予測データ毎に評価値を算出する。例えば、判定部 254 は、予測データ毎にクラスタ間距離およびクラスタ内距離を算出し、クラスタ間距離およびクラスタ内距離を予測データの評価値とする。

30

【0090】

最終クラスタ計算部 255 は、分析対象データ 241 の最終的なクラスタ分析結果となる最終データ 246 を生成する処理部である。最終クラスタ計算部 255 は、評価値データテーブル 244 から中間データテーブル 245 を生成する処理を行った後に、中間データテーブル 245 を基にして、最終データ 246 a , 246 b , 246 c を生成する。

40

【0091】

最終クラスタ計算部 255 が、評価値データテーブル 244 から中間データテーブル 245 を生成する処理について説明する。最終クラスタ計算部 255 は、評価データテーブル 244 の予測データ毎の評価値を比較し、パレート解となる予測データを特定し、特定したパレート解となる予測データを、中間データテーブル 245 に設定する。

【0092】

図 12 は、各予測データのクラスタ内距離とクラスタ間距離との関係を示す図 (2) である。図 12 において、縦軸はクラスタ内距離を示し、横軸はクラスタ間距離を示す。一般的に、クラスタ間距離が大きいほど、また、クラスタ内距離が小さいほど、中間データは、良い予測データであると言える。このため、図 12 に示す例では、最終クラスタ計算

50

部 2 5 5 は、中間データ 2 4 3 a , 2 4 3 c , 2 4 3 f , 2 4 3 g , 2 4 3 z を、パレート解として特定する。

【 0 0 9 3 】

続いて、最終クラスタ計算部 2 5 5 が、中間データテーブル 2 4 5 から最終データ 2 4 6 を生成する処理について説明する。最終クラスタ計算部 2 5 5 は、中間データテーブル 2 4 5 の各予測データの評価値を比較して、類似する予測データ同士を同一グループに分類する処理を行う。例えば、最終クラスタ計算部 2 5 5 は、各予測データのクラスタ間距離の差分が閾値未満となり、かつ、各予測データのクラスタ内距離の差分が閾値未満となる予測データを、同一のグループに分類する。

【 0 0 9 4 】

図 1 2 に示す例では、最終クラスタ計算部 2 5 5 は、予測データ 2 4 3 a , 2 4 3 c をグループ 5 0 a に分類し、予測データ 2 4 3 f , 2 4 3 g をグループ 5 0 b に分類し、予測データ 2 4 3 i , 2 4 3 z をグループ 5 0 c に分類する。最終クラスタ計算部 2 5 5 は、分類したグループ毎に、最終データ 2 4 6 を生成する。

【 0 0 9 5 】

例えば、最終クラスタ計算部 2 5 5 は、グループ 5 0 a に含まれる予測データ 2 4 3 a , 2 4 3 c を基にして、最終データ 2 4 6 a を生成する。最終クラスタ計算部 2 5 5 は、グループ 5 0 b に含まれる予測データ 2 4 3 f , 2 4 3 g を基にして、最終データ 2 4 6 b を生成する。最終クラスタ計算部 2 5 5 は、グループ 5 0 c に含まれる予測データ 2 4 3 i , 2 4 3 z を基にして、最終データ 2 4 6 c を生成する。

【 0 0 9 6 】

最終クラスタ計算部 2 5 5 が、予測データを基にして、最終データを特定する処理は、図 1 の最終クラスタ計算部 1 5 5 が、中間データテーブル 1 4 5 の予測データを基にして、最終データを特定する処理と同様である。図 1 2 に示す例では、グループ 5 0 a , 5 0 b , 5 0 c について、最終データが特定され、最終データ 2 4 6 a , 2 4 6 b , 2 4 6 c が生成される。

【 0 0 9 7 】

次に、本実施例 2 に係る分析装置 2 0 0 の処理手順について説明する。図 1 3 および図 1 4 は、実施例 2 にかかる分析装置の処理手順を示すフローチャートである。図 1 3 に示すように、分析装置 2 0 0 は、分析対象データ 2 4 1 を受け付ける（ステップ S 2 0 1）。また、分析装置 2 0 0 は、繰り返し計算回数を受け付け（ステップ S 2 0 2）、サンプリング件数を受け付ける。また、カウント値を初期化する（ステップ S 2 0 3）。分析装置 2 0 0 は、カウント値に 1 を加算する（ステップ S 2 0 4）。カウント値の初期値を 0 とする。

【 0 0 9 8 】

分析装置 2 0 0 は、分析対象データ 2 4 1 をサンプリングし、サンプリングデータを生成する（ステップ S 2 0 5）。各サンプリングデータは、サンプリングデータテーブル 2 4 2 に格納される。分析装置 2 0 0 は、サンプリングデータに対してクラスタ分析処理を実行し、各々のレコードに対してクラスタ番号を割り振る（ステップ S 2 0 6）。

【 0 0 9 9 】

分析装置 2 0 0 は、クラスタ番号を割り振ったサンプリングデータと分析対象データ 2 4 1 とを比較して、分析対象データ 2 4 1 に含まれる各々のレコードに対してクラスタ番号を割り振ることで予測データを生成する（ステップ S 2 0 7）。各予測データは、予測データテーブル 2 4 3 に格納される。

【 0 1 0 0 】

分析装置 2 0 0 は、予測データを基にして、クラスタ内距離およびクラスタ間距離を算出し、評価値データテーブル 2 4 4 を生成する（ステップ S 2 0 8）。分析装置 2 0 0 は、繰り返しの計算回数がカウント値未満であるか否かを判定する（ステップ S 2 0 9）。分析装置 2 1 0 0 は、繰り返しの計算回数がカウント値未満の場合には（ステップ S 2 0 9 0 , Y e s ）、ステップ S 2 0 4 に移行する。

10

20

30

40

50

【0101】

一方、分析装置200は、繰り返しの計算回数がカウント値以上である場合には（ステップS209，No）、図14のステップS210に移行する。

【0102】

図14の説明に移行する。分析装置200は、パレート解に対応する予測データを選択して、中間データテーブル245を作成する（ステップS210）。分析装置200は、パレート解に対応する各予測データの類似度を算出する（ステップS211）。分析装置200は、類似する各予測データを、グループに分類する（ステップS212）。

【0103】

分析装置200は、未選択のグループを選択し（ステップS213）、ランダムにクラスタ番号を割り振った複数の最終データ候補を生成する（ステップS214）。ステップS214において、分析装置200は、類似度が大きくなるようにクラスタ番号を割り振る。例えば、分析装置200は、ランダムにクラスタ番号を割り振り、類似度を計算する。そして、分析装置200は、類似度が大きい、クラスタ番号の割り振りを少し変更して、類似度が大きくなるか、試行する処理を利用者が設定した回数繰り返す。

10

【0104】

分析装置200は、グループに含まれる予測データと各最終データ候補とを比較して、類似度が最大となる最終データ候補を判定する（ステップS215）。

【0105】

分析装置200は、未選択のグループが存在するか否かを判定する（ステップS216）。分析装置200は、未選択のグループが存在する場合には（ステップS216，Yes）、ステップS213に移行する。一方、分析装置200は、未選択のグループが存在しない場合には（ステップS216，No）、各グループの判定結果を出力する（ステップS217）。

20

【0106】

次に、本実施例2に係る分析装置200の効果について説明する。分析装置200は、複数の予測データのうち、評価値のよい予測データを類似する予測データ同士でグルーピングし、グループ毎に、最終データ246を生成する。このため、分析装置200によれば、クラスタ分析に要する時間を削減することができる。また、類似する予測データに応じた最終データの候補を複数得ることが出来る。

30

【0107】

次に、上記実施例に示した分析装置100，200と同様の機能を実現する分析プログラムを実行するコンピュータの一例について説明する。図15は、分析プログラムを実行するコンピュータの一例を示す図である。

【0108】

図15に示すように、コンピュータ300は、各種演算処理を実行するCPU301と、ユーザからのデータの入力を受け付ける入力装置302と、ディスプレイ303を有する。また、コンピュータ300は、記憶媒体からプログラム等を読み取る読み取り装置304と、ネットワークを介して他のコンピュータとの間でデータの授受を行うインターフェース装置305とを有する。また、コンピュータ300は、各種情報を一時記憶するRAM306と、ハードディスク装置307を有する。そして、各装置301～307は、バス308に接続される。

40

【0109】

ハードディスク装置307は、サンプリングプログラム307a、クラスタ分析プログラム307b、クラスタ予測プログラム307c、判定プログラム307d、最終クラスタ計算プログラム307eを有する。CPU301は、各プログラム307a～307eを読み出してRAM306に展開する。

【0110】

サンプリングプログラム307aは、サンプリングプロセス306aとして機能する。クラスタ分析プログラム307bは、クラスタ分析プロセス306bとして機能する。ク

50

ラスト予測プログラム 307c は、クラスタ予測プロセス 306c として機能する。判定プログラム 307d は、判定プロセス 306d として機能する。最終クラスタ計算プログラム 307e は、最終クラスタ計算プロセス 306e として機能する。

【0111】

例えば、サンプリングプロセス 306a は、サンプリング実行部 151, 251 に対応する。クラスタ分析プロセス 306b は、クラスタ分析部 152, 252 に対応する。クラスタ予測プロセス 306c は、クラスタ予測部 153, 253 に対応する。判定プロセス 306d は、判定部 154, 254 に対応する。最終クラスタ計算プロセス 306e は、最終クラスタ計算部 155, 255 に対応する。

【0112】

なお、各プログラム 307a ~ 307e については、必ずしも最初からハードディスク装置 307 に記憶させておかなくても良い。例えば、コンピュータ 300 に挿入されるフレキシブルディスク (FD)、CD-ROM、DVD ディスク、光磁気ディスク、IC カードなどの「可搬用の物理媒体」に各プログラムを記憶させておく。そして、コンピュータ 500 がこれらから各プログラム 307a ~ 307e を読み出して実行するようにしてもよい。

【0113】

以上の各実施例を含む実施形態に関し、さらに以下の付記を開示する。

【0114】

(付記 1) 入力データに対してサンプリングを実行し前記入力データから一部のデータを抽出する処理を繰り返し実行して複数のサンプリングデータを生成するサンプリング実行部と、

前記複数のサンプリングデータについてクラスタ分析を実行し、前記サンプリングデータ毎に、前記サンプリングデータに含まれるデータを異なるクラスタに分類するクラスタ分析部と、

前記複数のサンプリングデータに対する前記クラスタ分析部の複数の分類結果と前記入力データとを基にして、前記入力データに含まれるデータの所属するクラスタを予測したデータを示す予測データを複数生成するクラスタ予測部と、

前記予測データのクラスタ間距離およびクラスタ内距離を基にして、予測データ毎に評価値を算出し、パレート解となる評価値に対応する予測データを判定する判定部と、

前記パレート解となる評価値に対応する予測データを基にして、前記入力データに含まれるデータをクラスタに分類する最終クラスタ計算部と

を有することを特徴とする分析装置。

【0115】

(付記 2) 前記最終クラスタ計算部は、パレート解となる評価値に対応する類似の予測データをグループ化し、同一グループに含まれる予測データを基にして、前記入力データに含まれるデータを異なるクラスタに分類する処理を、グループ毎に実行することを特徴とする付記 1 に記載の分析装置。

【0116】

(付記 3) 前記最終クラスタ計算部は、前記入力データに対して、ランダムにクラスタを割り当てた複数の最終クラスタデータを生成し、各最終クラスタデータと予測データとの類似度を基にして、特定の最終クラスタデータを選択することを特徴とする付記 1 または 2 に記載の分析装置。

【0117】

(付記 4) コンピュータが実行する分析方法であって、

入力データに対してサンプリングを実行し前記入力データから一部のデータを抽出する処理を繰り返し実行して複数のサンプリングデータを生成し、

前記複数のサンプリングデータについてクラスタ分析を実行し、前記サンプリングデータ毎に、前記サンプリングデータに含まれるデータを異なるクラスタに分類し、

前記複数のサンプリングデータに対する前記クラスタ分析部の複数の分類結果と前記入

10

20

30

40

50

力データとを基にして、前記入力データに含まれるデータの所属するクラスタを予測したデータを示す予測データを複数生成し、

前記予測データのクラスタ間距離およびクラスタ内距離を基にして、予測データ毎に評価値を算出し、パレート解となる評価値に対応する予測データを判定し、

前記パレート解となる評価値に対応する予測データを基にして、前記入力データに含まれるデータをクラスタに分類する

各処理を実行することを特徴とする分析方法。

【0118】

(付記5) 前記入力データに含まれるデータをクラスタに分類する処理は、パレート解となる評価値に対応する類似の予測データをグループ化し、同一グループに含まれる予測データを基にして、前記入力データに含まれるデータを異なるクラスタに分類する処理を、グループ毎に実行することを特徴とする付記4に記載の分析方法。

10

【0119】

(付記6) 前記入力データに含まれるデータをクラスタに分類する処理は、前記入力データに対して、ランダムにクラスタを割り当てた複数の最終クラスタデータを生成し、各最終クラスタデータと予測データとの類似度を基にして、特定の最終クラスタデータを選択することを特徴とする付記4または5に記載の分析方法。

【0120】

(付記7) コンピュータに、

入力データに対してサンプリングを実行し前記入力データから一部のデータを抽出する処理を繰り返し実行して複数のサンプリングデータを生成し、

20

前記複数のサンプリングデータについてクラスタ分析を実行し、前記サンプリングデータ毎に、前記サンプリングデータに含まれるデータを異なるクラスタに分類し、

前記複数のサンプリングデータに対する前記クラスタ分析部の複数の分類結果と前記入力データとを基にして、前記入力データに含まれるデータの所属するクラスタを予測したデータを示す予測データを複数生成し、

前記予測データのクラスタ間距離およびクラスタ内距離を基にして、予測データ毎に評価値を算出し、パレート解となる評価値に対応する予測データを判定し、

前記パレート解となる評価値に対応する予測データを基にして、前記入力データに含まれるデータをクラスタに分類する

30

各処理を実行させることを特徴とする分析プログラム。

【0121】

(付記8) 前記入力データに含まれるデータをクラスタに分類する処理は、パレート解となる評価値に対応する類似の予測データをグループ化し、同一グループに含まれる予測データを基にして、前記入力データに含まれるデータを異なるクラスタに分類する処理を、グループ毎に実行することを特徴とする付記7に記載の分析プログラム。

【0122】

(付記9) 前記入力データに含まれるデータをクラスタに分類する処理は、前記入力データに対して、ランダムにクラスタを割り当てた複数の最終クラスタデータを生成し、各最終クラスタデータと予測データとの類似度を基にして、特定の最終クラスタデータを選択することを特徴とする付記4または5に記載の分析方法。

40

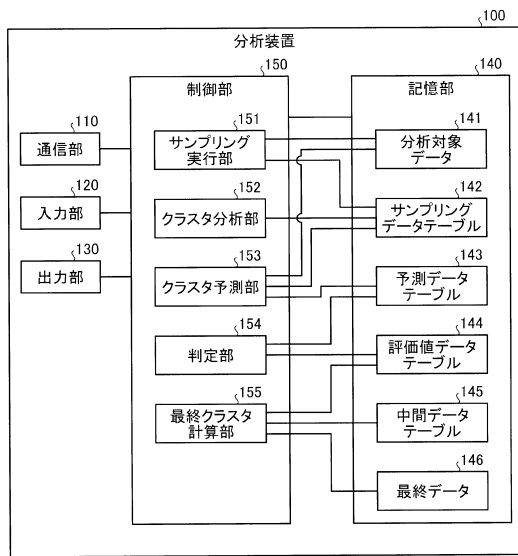
【符号の説明】

【0123】

- 100, 200 分析装置
- 151, 251 サンプリング実行部
- 152, 252 クラスタ分析部
- 153, 253 クラスタ予測部
- 154, 254 判定部
- 155, 255 最終クラスタ計算部

【図 1】

本実施例1に係る分析装置の構成を示す機能ブロック図



【図 2】

分析対象データのデータ構造の一例を示す図

識別番号	年齢	性別	身長	体重	...
1001	35	1	170	65	...
1002	40	2	165	59	...
...					

【図 3】

サンプリングデータテーブルのデータ構造の一例を示す図

識別番号	年齢	性別	身長	体重	...	クラス番号
1001	35	1	170	65	...	1
1002	40	2	180	70	...	2
...						

【図 4】

予測データテーブルのデータ構造の一例を示す図

識別番号	年齢	性別	身長	体重	...	クラス番号
1001	35	1	170	65	...	1
1002	40	2	165	59	...	1
...						
1005	44	2	180	70	...	2
...						

【図 6】

中間データテーブルのデータ構造の一例を示す図

識別番号	クラス番号:1	クラス番号:2	クラス番号:3
1001	1	0	0
1002	1	0	0
1003	0	1	0
1004	0	0	1
...			

【図 5】

評価値データテーブルのデータ構造の一例を示す図

予測データ 識別情報	評価値	
	クラス間距離	クラス内距離
予測データ1	5	25
予測データ2	10	20
予測データ3	16	15
...		

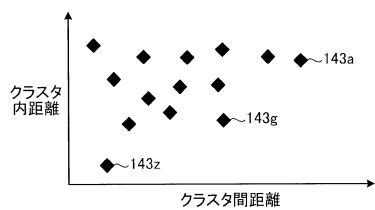
【図 7】

最終データのデータ構造の一例を示す図

識別番号	クラス番号:1	クラス番号:2	クラス番号:3
1001	1	0	0
1002	1	0	0
1003	0	1	0
1004	0	0	1
...			

【図 8】

各予測データのクラスタ内距離とクラスタ間距離との関係を示す図(1)



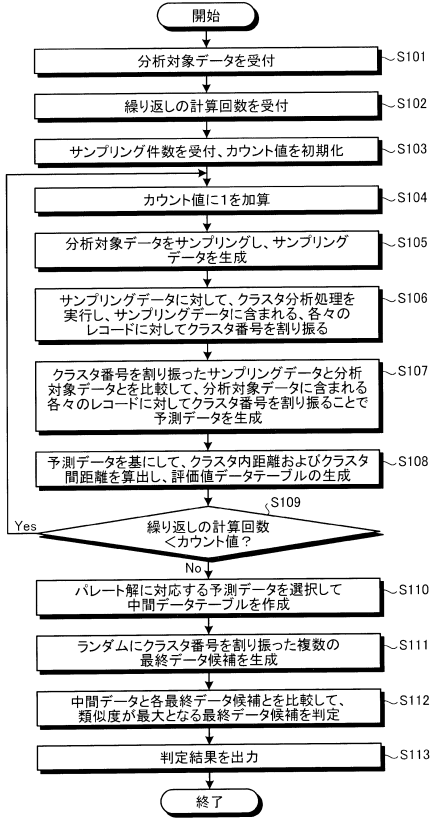
【図 9】

最終データ候補テーブルの一例を示す図

10			
10a			
識別番号	クラスタ番号:1	クラスタ番号:2	クラスタ番号:3
1001	1	0	0
1002	1	0	0
1003	0	1	0
1004	0	0	1
...			
10b			
識別番号	クラスタ番号:1	クラスタ番号:2	クラスタ番号:3
1001	0	0	1
1002	1	0	0
1003	0	1	0
1004	0	0	1
...			
10c			
識別番号	クラスタ番号:1	クラスタ番号:2	クラスタ番号:3
1001	0	1	0
1002	0	1	0
1003	0	1	0
1004	0	0	1
...			
:			

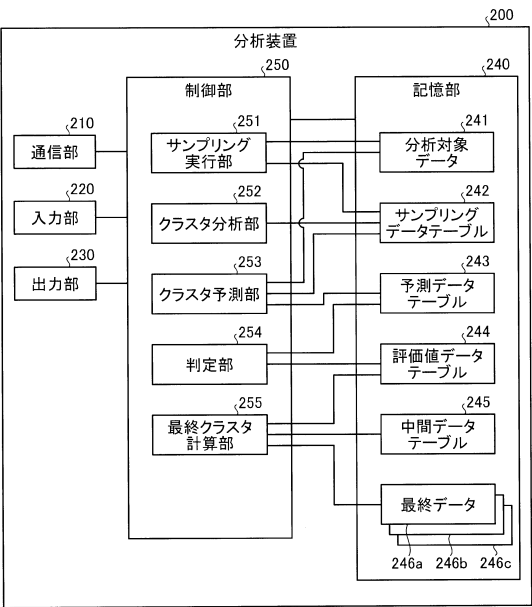
【図 10】

本実施例1にかかる分析装置の処理手順を示すフローチャート



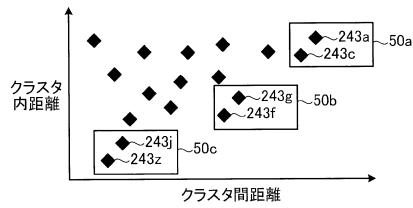
【図 11】

本実施例2にかかる分析装置の構成を示す機能ブロック図



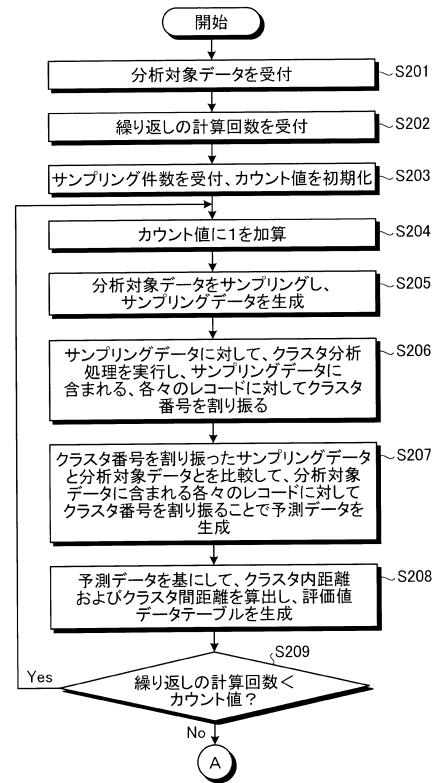
【図 12】

各予測データのクラスタ内距離とクラスタ間距離との関係を示す図(2)



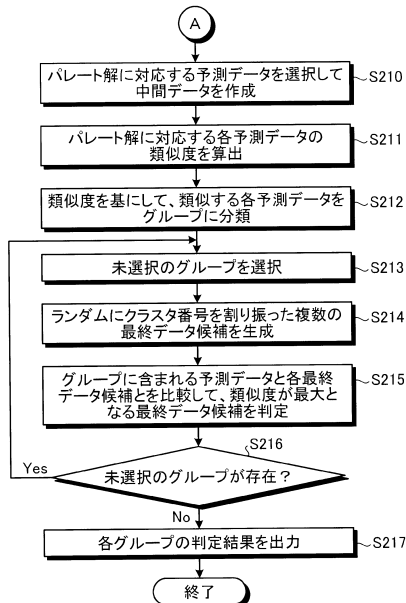
【図 13】

実施例2にかかる分析装置の処理手順を示すフローチャート(1)



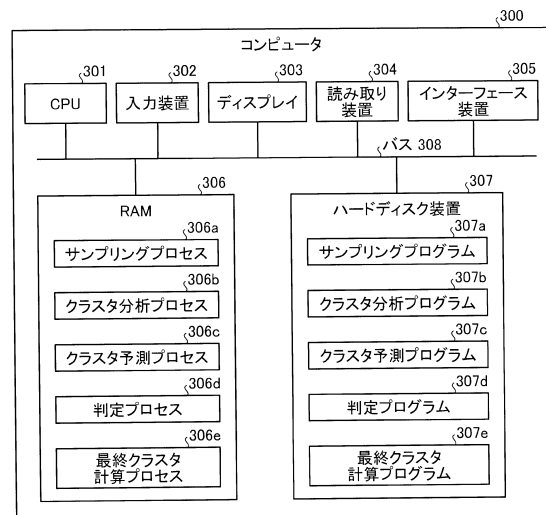
【図 14】

実施例2にかかる分析装置の処理手順を示すフローチャート(2)



【図 15】

分析プログラムを実行するコンピュータの一例を示す図



フロントページの続き

(56)参考文献 特開2007-179143(JP,A)

特開平11-250094(JP,A)

米国特許第7162413(US,B1)

特開2007-179288(JP,A)

特開2007-172427(JP,A)

小林 優、外2名、特徴要素の重みを考慮に入れたクラスタ代表の洗練による文書クラスタリング、情報処理学会研究報告、日本、社団法人情報処理学会、2002年 3月15日、第2002巻、第28号、p.135-142

山代 大輔、外2名、可視化手法を用いた多目的最適化問題における満足解の選択支援、知能と情報、日本、日本知能情報ファジィ学会、2008年12月15日、第20巻、第6号、p.24-32

(58)調査した分野(Int.Cl., DB名)

G06F 17/30