

(19)日本国特許庁(JP)

(12)特許公報(B2)

(11)特許番号
特許第7586196号
(P7586196)

(45)発行日 令和6年11月19日(2024.11.19)

(24)登録日 令和6年11月11日(2024.11.11)

(51)国際特許分類 F I
G 0 6 Q 10/063 (2023.01) G 0 6 Q 10/063

請求項の数 10 (全24頁)

(21)出願番号	特願2022-571910(P2022-571910)	(73)特許権者	000004237 日本電気株式会社 東京都港区芝五丁目7番1号
(86)(22)出願日	令和3年10月25日(2021.10.25)	(74)代理人	110000338 弁理士法人 H A R A K E N Z O W O R L D P A T E N T & T R A D E M A R K
(86)国際出願番号	PCT/JP2021/039367	(72)発明者	野澤 拓磨 東京都港区芝五丁目7番1号 日本電気 株式会社内
(87)国際公開番号	WO2022/137778	(72)発明者	小山田 昌史 東京都港区芝五丁目7番1号 日本電気 株式会社内
(87)国際公開日	令和4年6月30日(2022.6.30)	(72)発明者	董 于洋 東京都港区芝五丁目7番1号 日本電気 株式会社内
審査請求日	令和5年6月14日(2023.6.14)		
(31)優先権主張番号	特願2020-212788(P2020-212788)		
(32)優先日	令和2年12月22日(2020.12.22)		
(33)優先権主張国・地域又は機関	日本国(JP)		

最終頁に続く

(54)【発明の名称】 情報処理装置、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】

複数のデータセットのそれぞれから当該データセットに含まれる複数のデータ項目を関連付けることにより生成されたデータであるインサイトサブジェクトを、検出対象のインサイトごとにグループ化する分類手段と、

グループ化された複数の前記インサイトサブジェクトの組み合わせについて、インサイトの有無を判定するための評価値を算出する評価手段と、
複数の前記インサイトサブジェクトにおけるデータの粒度を統一する粒度統一手段と、を備え、

前記評価手段は、粒度が統一された複数の前記インサイトサブジェクトについて前記評価値を算出する、
情報処理装置。

【請求項2】

複数の前記インサイトサブジェクトにおける表記を統一する表記統一手段を備え、
前記分類手段は、表記が統一された前記インサイトサブジェクトをグループ化する、請求項1に記載の情報処理装置。

【請求項3】

前記評価手段は、動的時間伸縮法または関数データ解析により前記評価値を算出する、請求項1に記載の情報処理装置。

【請求項4】

前記評価手段は、グループ化された複数の前記インサイトサブジェクトを主成分分析することにより求めた、各主成分の寄与度の偏りの程度に基づいて前記評価値を算出する、請求項 1 から 3 の何れか 1 項に記載の情報処理装置。

【請求項 5】

前記主成分分析により求められた主成分を用いて、グループ化された複数の前記インサイトサブジェクトに含まれるデータを表すことにより、当該データに含まれる外れ値を検出する外れ値検出手段を備える、請求項 4 に記載の情報処理装置。

【請求項 6】

少なくとも 1 つのプロセッサが、
複数のデータセットのそれぞれから当該データセットに含まれる複数のデータ項目を関連付けることにより生成されたデータであるインサイトサブジェクトを、検出対象のインサイトごとにグループ化すること、
グループ化された複数の前記インサイトサブジェクトの組み合わせについて、インサイトの有無を判定するための評価値を算出すること、および
複数の前記インサイトサブジェクトにおけるデータの粒度を統一すること、を含み、
前記評価値を算出する工程において、前記少なくとも 1 つのプロセッサが、粒度が統一された複数の前記インサイトサブジェクトについて前記評価値を算出する、
分析方法。

【請求項 7】

コンピュータに、
複数のデータセットのそれぞれから当該データセットに含まれる複数のデータ項目を関連付けることにより生成されたデータであるインサイトサブジェクトを、検出対象のインサイトごとにグループ化する処理と、
グループ化された複数の前記インサイトサブジェクトの組み合わせについて、インサイトの有無を判定するための評価値を算出する処理と、
複数の前記インサイトサブジェクトにおけるデータの粒度を統一する処理と、を実行させる分析プログラムであって、
前記評価値を算出する処理において、前記コンピュータは、粒度が統一された複数の前記インサイトサブジェクトについて前記評価値を算出する、
分析プログラム。

【請求項 8】

複数のデータセットのそれぞれから当該データセットに含まれる複数のデータ項目を関連付けることにより生成されたデータであるインサイトサブジェクトを、検出対象のインサイトごとにグループ化する分類手段と、
グループ化された複数の前記インサイトサブジェクトの組み合わせについて、インサイトの有無を判定するための評価値を算出する評価手段であって、グループ化された複数の前記インサイトサブジェクトを主成分分析することにより求めた、各主成分の寄与度の偏りの程度に基づいて前記評価値を算出する評価手段と、
前記主成分分析により求められた主成分を用いて、グループ化された複数の前記インサイトサブジェクトに含まれるデータを表すことにより、当該データに含まれる外れ値を検出する外れ値検出手段と、を備える情報処理装置。

【請求項 9】

少なくとも 1 つのプロセッサが、
複数のデータセットのそれぞれから当該データセットに含まれる複数のデータ項目を関連付けることにより生成されたデータであるインサイトサブジェクトを、検出対象のインサイトごとにグループ化すること、および
グループ化された複数の前記インサイトサブジェクトの組み合わせについて、インサイトの有無を判定するための評価値を算出すること、を含み、
前記評価値を算出する工程において、前記少なくとも 1 つのプロセッサが、グループ化された複数の前記インサイトサブジェクトを主成分分析することにより求めた、各主成分の

10

20

30

40

50

寄与度の偏りの程度に基づいて前記評価値を算出し、
 前記少なくとも1つのプロセッサが、前記主成分分析により求められた主成分を用いて、
 グループ化された複数の前記インサイトサブジェクトに含まれるデータを表すことにより、
 当該データに含まれる外れ値を検出すること、
 を含む分析方法。

【請求項10】

コンピュータに、

複数のデータセットのそれぞれから当該データセットに含まれる複数のデータ項目を関連
 付けることにより生成されたデータであるインサイトサブジェクトを、検出対象のインサ
 イトごとにグループ化する処理と、

10

グループ化された複数の前記インサイトサブジェクトの組み合わせについて、インサイト
 の有無を判定するための評価値を算出する処理であって、グループ化された複数の前記イン
 サイトサブジェクトを主成分分析することにより求めた、各主成分の寄与度の偏りの程
 度に基づいて前記評価値を算出する処理と、

前記主成分分析により求められた主成分を用いて、グループ化された複数の前記インサ
 イトサブジェクトに含まれるデータを表すことにより、当該データに含まれる外れ値を検出
 する処理と、を実行させる分析プログラム。

【発明の詳細な説明】

【技術分野】

【0001】

本発明は、データセットの解析を行う情報処理装置等に関する。

20

【背景技術】

【0002】

近年、様々な分野において、データを収集し、そのデータを分析することにより、人に
 にとって意味のある知見を見出すことが行われている。このような知見はインサイトと呼ば
 れる。一般的なデータ分析作業では、分析者が、仮説を設定し、設定した仮説に基づいて
 データ分析・可視化し、その仮説検証する、というサイクルを繰り返すことによってイン
 サイトを見出している。

【0003】

インサイトを見出すための上記のようなデータ分析作業は、非常に時間と労力を要する
 ものであるため、これを自動化する技術の開発が進められている。例えば、下記の特許文
 献1には、データセットから自動でインサイトを提供するシステムが開示されている。分
 析者は、特許文献1に記載のシステムに、分析したい多次元データを入力すればよい。こ
 れにより、当該システムにより自動的にインサイトが決定され、決定されたインサイトが
 ディスプレイに表示される。

30

【先行技術文献】

【特許文献】

【0004】

【文献】米国特許第2020/0257682号明細書

【発明の概要】

40

【発明が解決しようとする課題】

【0005】

特許文献1に記載の技術には、複数のデータセット間のインサイトを検出することがで
 きないという点で改善の余地があった。例えば、ある企業の製品販売データからなるデー
 タセットと、他の企業についての製品販売データからなるデータセットの両方を解析する
 ことにより、一方のデータセットのみからは得られないインサイトが見つかる可能性があ
 る。

【0006】

しかしながら、特許文献1に記載の技術では、このような複数のデータセット間のイン
 サイトを検出することは想定されていない。このため、当然のことながら、特許文献1に

50

記載の技術では、複数のデータセット間のインサイトを検出することはできない。

【0007】

本発明の一態様は、上記の問題に鑑みてなされたものであり、その目的の一例は、複数のデータセット間におけるインサイトの検出を可能にする情報処理装置等を提供することである。

【課題を解決するための手段】

【0008】

本発明の一態様に係る情報処理装置は、複数のデータセットのそれぞれから当該データセットに含まれる複数のデータ項目を関連付けることにより生成されたデータであるインサイトサブジェクトを、検出対象のインサイトごとにグループ化する分類手段と、グループ化された複数の前記インサイトサブジェクトの組み合わせについて、インサイトの有無を判定するための評価値を算出する評価手段とを備える。

10

【0009】

本発明の一態様に係る分析方法は、少なくとも1つのプロセッサが、複数のデータセットのそれぞれから当該データセットに含まれる複数のデータ項目を関連付けることにより生成されたデータであるインサイトサブジェクトを、検出対象のインサイトごとにグループ化することと、グループ化された複数の前記インサイトサブジェクトの組み合わせについて、インサイトの有無を判定するための評価値を算出すること、を含む。

【0010】

本発明の一態様に係る分析プログラムは、複数のデータセットのそれぞれから当該データセットに含まれる複数のデータ項目を関連付けることにより生成されたデータであるインサイトサブジェクトを、検出対象のインサイトごとにグループ化する処理と、グループ化された複数の前記インサイトサブジェクトの組み合わせについて、インサイトの有無を判定するための評価値を算出する処理と、をコンピュータに実行させる。

20

【発明の効果】

【0011】

本発明の一態様によれば、複数のデータセット間におけるインサイトの検出が可能になる。

【図面の簡単な説明】

【0012】

【図1】本発明の例示的实施形態1に係る情報処理装置の構成を示すブロック図である。

【図2】本発明の例示的实施形態1に係る分析方法の流れを示すフロー図である。

【図3】本発明の例示的实施形態2に係る情報処理装置が実行する処理の概要を示す図である。

30

【図4】本発明の例示的实施形態2に係る情報処理装置の構成を示すブロック図である。

【図5】本発明の例示的实施形態2に係る分析方法の流れを示すフロー図である。

【図6】分析対象データと、当該分析対象データから生成されたインサイトサブジェクトの例を示す図である。

【図7】評価結果データと出力データの例を示す図である。

【図8】本発明の例示的实施形態3に係る情報処理装置の構成を示すブロック図である。

40

【図9】本発明の例示的实施形態3に係る分析方法の流れを示すフロー図である。

【図10】インサイトスコアの算出方法と、外れ値の検出方法を説明する図である。

【図11】上記情報処理装置の各機能を実現するソフトウェアであるプログラムの命令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】

【0013】

〔例示的实施形態1〕

本発明の第1の例示的实施形態について、図面を参照して詳細に説明する。本例示的实施形態は、後述する例示的实施形態の基本となる形態である。

【0014】

50

(情報処理装置 1 の構成)

本例示的实施形態に係る情報処理装置 1 の構成について、図 1 を参照して説明する。図 1 は、情報処理装置 1 の構成を示すブロック図である。図示のように、情報処理装置 1 は、分類部 1 1 と評価部 1 2 を備えている。

【 0 0 1 5 】

分類部 1 1 は、複数のデータセットのそれぞれから当該データセットに含まれる複数のデータ項目を関連付けることにより生成されたデータであるインサイトサブジェクトを、検出対象のインサイトごとにグループ化する。グループ化の際に、分類部 1 1 は、評価部 1 2 による評価値の算出が可能なインサイトサブジェクトをグループ化する。なお、以下では、検出対象のインサイトをインサイトタイプと呼ぶ。インサイトタイプは少なくとも 1 つ設定されていけばよい。インサイトタイプの詳細は例示的实施形態 2 で説明する。

10

【 0 0 1 6 】

そして、評価部 1 2 は、グループ化された複数の前記インサイトサブジェクトの組み合わせについて、インサイトの有無を判定するための評価値を算出する。以下では、この評価値をインサイトスコアと呼ぶ。

【 0 0 1 7 】

例えば、ある店舗の月間の売上記録を示すデータセットが分析対象である場合、その店舗における日別の総売上を示すデータ（日付と総売上のデータ項目を関連付けたデータ）をインサイトサブジェクトとすることができる。同様に、その店舗におけるある商品の日別の売上を示すデータ（日付とある商品の売上のデータ項目を関連付けたデータ）をインサイトサブジェクトとすることができる。このようなインサイトサブジェクトは、例えばチャート等の形式で可視化することができるため、インサイトサブジェクトを可視化パターンと呼ぶこともできる。インサイトサブジェクトは、多次元データであるデータセットから得られる各可視化パターンを特徴づけるものであるということもできる。この場合、1 つのインサイトサブジェクトにつき 1 つの可視化パターンが対応付けられる。

20

【 0 0 1 8 】

そして、検出対象のインサイト、すなわちインサイトタイプが、例えばインサイトサブジェクト間の相関であれば、分類部 1 1 は、相関の有無を判定するためのインサイトスコア（例えば相関係数）の算出が可能なインサイトサブジェクトをグループ化する。例えば、分類部 1 1 は、上記の例では、各店舗における日付と売上の関係を示すインサイトサブジェクトをグループ化してもよい。これにより、評価部 1 2 は、各店舗における日付と売上についてインサイトスコアを算出することができる。インサイトスコアは、そのまま出力してもユーザがインサイトを発見する大きな助けとなる。また、インサイトスコアを用いることにより、インサイトスコアが高い、すなわちインサイトである可能性が高いインサイトサブジェクトの組み合わせを自動で検出することも可能になる。

30

【 0 0 1 9 】

以上のように、本例示的实施形態に係る情報処理装置 1 では、複数のデータセットのそれぞれから生成されたインサイトサブジェクトを、検出対象のインサイトごとにグループ化する分類部 1 1 と、グループ化された複数の前記インサイトサブジェクトの組み合わせについて、インサイトの有無を判定するための評価値を算出する評価部 1 2 と、を備える、という構成が採用されている。

40

【 0 0 2 0 】

したがって、本例示的实施形態に係る情報処理装置 1 によれば、複数のデータセット間におけるインサイトの検出が可能になるという効果が得られる。言い換えれば、本例示的实施形態に係る情報処理装置 1 によれば、複数のデータセットを横断的に分析することで得られる複合インサイト（以下、横断的複合インサイトと呼ぶ）の発見に繋がる可能性のあるデータをユーザに提示することが可能になる。

【 0 0 2 1 】

なお、上述の情報処理装置 1 の機能は、プログラムによって実現することもできる。本例示的实施形態に係る分析プログラムは、コンピュータに、複数のデータセットのそれぞ

50

れから生成されたインサイトサブジェクトを、検出対象のインサイトごとにグループ化する処理と、グループ化された複数の前記インサイトサブジェクトの組み合わせについて、インサイトの有無を判定するための評価値を算出する処理と、を実行させる。したがって、本例示的实施形態に係る分析プログラムによれば、複数のデータセット間におけるインサイト、すなわち横断的複合インサイトの検出が可能になるという効果が得られる。

【0022】

(分析方法の流れ)

本例示的实施形態に係る分析方法の流れについて、図2を参照して説明する。図2は、本例示的实施形態に係る分析方法の流れを示すフロー図である。

【0023】

S11では、少なくとも1つのプロセッサが、複数のデータセットのそれぞれから生成されたインサイトサブジェクトを、インサイトタイプごとにグループ化する。そして、S12では、少なくとも1つのプロセッサが、S11でグループ化された複数の前記インサイトサブジェクトの組み合わせについて、インサイトの有無を判定するための評価値であるインサイトスコアを算出する。これにより、図2の分析方法は終了する。

【0024】

なお、1つのプロセッサにS11～S12の処理を実行させてもよいし、S11の処理とS12の処理をそれぞれ別のプロセッサに実行させてもよい。後者の場合、各プロセッサは、1つの情報処理装置が備えているものであってもよいし、それぞれ異なる情報処理装置が備えているものであってもよい。また、S11～S12の処理を実行する少なくとも1つのプロセッサは、情報処理装置1が備えているものであってもよい。

【0025】

以上のように、本例示的实施形態に係る分析方法においては、少なくとも1つのプロセッサが、複数のデータセットのそれぞれから生成されたインサイトサブジェクトをインサイトタイプごとにグループ化すること、およびグループ化された複数の前記インサイトサブジェクトの組み合わせについて、インサイトの有無を判定するためのインサイトスコアを算出すること、を含む、という構成が採用されている。このため、本例示的实施形態に係る分析方法によれば、複数のデータセット間におけるインサイト、すなわち横断的複合インサイトの検出が可能になるという効果が得られる。

【0026】

(例示的实施形態2)

(概要)

本発明の第2の例示的实施形態について、図面を参照して詳細に説明する。本例示的实施形態では、複数のデータセットの入力を受け付けて、それらのデータセットについてのインサイトに関する情報を出力する情報処理装置2について説明する。図3は、情報処理装置2が実行する処理の概要を示す図である。

【0027】

まず、情報処理装置2は、分析対象となる分析対象データ211aと211bを取得する。分析対象データ211aと211bは、何れも複数のレコードを含む多次元データのデータセットである。なお、分析対象データ211aと211bを区別する必要がないときには単に分析対象データ211と記載する。図3に示す分析対象データ211aと211bは何れもテーブル形式のデータである。

【0028】

次に、情報処理装置2は、取得した分析対象データ211aと211bのそれぞれからインサイトサブジェクトを生成する。図3の例では、分析対象データ211aからI₁～I₃の3つのインサイトサブジェクトが生成され、分析対象データ211bからI₄、I₅の2つのインサイトサブジェクトが生成されている。

【0029】

続いて、情報処理装置2は、生成したインサイトサブジェクトI₁～I₅をグループ化する。図3の例では、インサイトサブジェクトI₁とI₅がグループG¹に分類され、イン

10

20

30

40

50

サイトサブジェクト I_3 と I_4 がグループ G^2 に分類されている。グループ G^1 と G^2 のインサイトタイプは同じであってもよいし、異なってもよい。ただし、グループ G^1 と G^2 のインサイトタイプが同じである場合には、各グループにはそれぞれ異なるインサイトサブジェクトを分類する。

【0030】

そして、情報処理装置2は、各グループに含まれるインサイトサブジェクトの組み合わせについて、インサイトの有無を判定するための評価値であるインサイトスコアを算出する。図3の例では、インサイトサブジェクト I_1 と I_5 のインサイトスコアが0.6、インサイトサブジェクト I_3 と I_4 のインサイトスコアが0.9と算出されている。インサイトスコアは、例えばインサイトサブジェクト間の相関の程度を0~1の数値(数値が大きいほど相関の程度が高い)で示すものであってもよい。この場合、インサイトサブジェクト I_3 と I_4 は、相関が高いことになる。

10

【0031】

ここで、インサイトサブジェクト I_3 は、分析対象データ211aから生成されたものである。一方、インサイトサブジェクト I_4 は、分析対象データ211bから生成されたものである。そして、インサイトサブジェクト I_3 と I_4 の相関が高いという知見は、人にとって有用なものである。つまり、情報処理装置2によれば、複数のデータセット間におけるインサイト、すなわち横断的複合インサイトの検出が可能になる。なお、詳細は以下説明するが、情報処理装置2は、相関以外にも様々なインサイトの検出を可能にする。

【0032】

20

(情報処理装置2の構成)

図4は、情報処理装置2の構成を示すブロック図である。情報処理装置2は、情報処理装置2の各部を統括して制御する制御部20と、情報処理装置2が使用する各種データを記憶する記憶部21を備えている。また、情報処理装置2は、情報処理装置2が他の装置と通信するための通信部22、情報処理装置2に対する入力を受け付ける入力部23、および情報処理装置2がデータを出力するための出力部24を備えている。以下では、出力部24がデータを表示出力する表示装置である例を説明するが、出力部24の出力態様は任意であり、例えば印字出力や音声出力等の態様でデータを出力するものであってもよい。また、入力部23と出力部24は、情報処理装置2に外付けされた、情報処理装置2の外部の機器であってもよい。

30

【0033】

制御部20には、データ取得部201、サブジェクト生成部202、表記統一部203、分類部204、粒度統一部205、評価部206、および出力データ生成部207が含まれている。また、記憶部21には、分析対象データ211、評価結果データ212、および出力データ213が記憶されている。

【0034】

分析対象データ211は、情報処理装置2による分析対象の対象となるデータである。分析対象データ211には、複数のデータセットが含まれている。各データセットは、複数のレコードを含む多次元データである。また、評価結果データ212は、評価部206による分析対象データ211の評価の結果を示すデータである。そして、出力データ213は、情報処理装置2による分析対象データ211の分析の結果をユーザに提示するためのデータ、すなわち分析対象データ211のインサイトに関するデータである。

40

【0035】

データ取得部201は、情報処理装置2が分析する対象となる複数のデータセットを取得し、それらを分析対象データ211として記憶部21に記憶させる。データ取得部201は、分析開始時までに分析対象データ211を取得して記憶部21に記憶させればよい。分析対象データ211の取得方法は特に限定されない。例えば、データ取得部201は、情報処理装置2のユーザが入力部23を介して入力したデータセットを取得してもよい。また、例えば、データ取得部201は、通信部22を介した通信により、外部の装置から分析対象データ211を取得してもよい。

50

【 0 0 3 6 】

サブジェクト生成部 2 0 2 は、分析対象データ 2 1 1 に含まれる複数のデータセットのそれぞれからインサイトサブジェクトを生成する。より詳細には、サブジェクト生成部 2 0 2 は、複数のデータセットのそれぞれから当該データセットに含まれる複数のデータ項目を関連付けることによりインサイトサブジェクトを生成する。例えば、あるデータセットが、日付、売上、および場所のデータ項目を含む多次元データである場合、サブジェクト生成部 2 0 2 は、日付と売上を関連付けたインサイトサブジェクトや、場所と売上を関連付けたインサイトサブジェクトを生成する。

【 0 0 3 7 】

表記統一部 2 0 3 は、各インサイトサブジェクトにおけるデータの表記を統一する。より詳細には、表記統一部 2 0 3 は、各インサイトサブジェクトに含まれる単語の中から類似した単語を抽出し、それらの単語を 1 つの単語に置き換えることにより、各インサイトサブジェクトにおける表記を統一する。なお、上記「類似」には、単語の文字列の類似の他、意味の類似も含まれる。

10

【 0 0 3 8 】

例えば、あるデータセットにおいて商品の販売地を表す「東京都」は、他のデータセットにおいて商品の販売地を表す「東京」と意味および文字列が類似した単語であり、これらは表記ゆれと呼ぶこともできる。また、例えば、あるデータセットにおいて商品の販売地を表す「都道府県」は、他のデータセットにおいて商品の販売地を表す「場所」と、意味が類似した単語である。

20

【 0 0 3 9 】

このような類似の単語を抽出する方法としては任意のものが適用可能である。表記統一部 2 0 3 は、「東京」と「東京都」のような表記ゆれの単語を抽出してもよい。この場合、表記統一部 2 0 3 は、例えば、単語間の編集距離が近い単語を抽出してもよい。編集距離は、レーベンシュタイン距離とも呼ばれ、2 つの文字列がどの程度異なっているかを示す距離である。編集距離を求める際には、表記統一部 2 0 3 は、比較対象の一方の単語を構成する文字列に対して何回の変更処理（削除、挿入、置換）を行えば、比較対象の他方を構成する文字列に変換できるかを求める。この他にも、分析対象データ 2 1 1 は、例えば 2 つの文字列の長さや置換の要不要（部分的な一致）を測る距離であるジャロ・ウィンクラー距離に基づいて類似の単語を抽出してもよい。

30

【 0 0 4 0 】

また、意味が類似した単語を抽出する場合、分析対象データ 2 1 1 は、例えば、各データセットに含まれる各単語を分散表現で表し、分散表現の類似度が高い単語を抽出してもよい。分散表現の導出には、例えば word2vec 等のプログラムを用いることができる。

【 0 0 4 1 】

表記統一部 2 0 3 は、類似した単語を抽出した後、それらの単語の表記を統一する。例えば、表記統一部 2 0 3 は、類似する 2 つの単語のうち一方の単語を他方の単語に全て置換することにより表記を統一してもよい。また、表記統一部 2 0 3 は、類似する 2 つの単語を、それらの単語を包括する上位概念的な単語に置換することにより表記を統一してもよい。

40

【 0 0 4 2 】

分類部 2 0 4 は、サブジェクト生成部 2 0 2 が生成したインサイトサブジェクトをグループ化する。より詳細には、分類部 2 0 4 は、インサイトの有無を判定するための評価値であるインサイトスコアを算出可能なインサイトサブジェクトをグループ化する。これにより、インサイトスコアに基づいてインサイトを検出することが可能になる。なお、1 つのグループには任意の数のインサイトサブジェクトを含めることができる。そして、1 つのグループには異なるデータセットから得られたインサイトサブジェクトを含めることができる。1 つのグループには少なくとも 1 つのインサイトサブジェクトを含めることが好ましい。

【 0 0 4 3 】

50

なお、表記統一部 203 が複数のインサイトサブジェクトにおける表記を統一していた場合、評価部 206 は、表記が統一されたインサイトサブジェクトをグループ化する。異なるデータセット間では、表記が不統一であることも多く、表記が不統一であることが評価の支障となることも一般的には多いが、情報処理装置 2 によればそのような場合にも評価を行うことができる。つまり、情報処理装置 2 によれば、例示的实施形態 1 に係る情報処理装置 1 の奏する効果に加えて、表記が不統一なデータセットについても横断的複合インサイトを検出することが可能になるという効果が得られる。

【0044】

例えば、年別の売上を示すインサイトサブジェクトが複数存在する場合、それらのインサイトサブジェクトの系列名は何れも「年」と「売上」となるから、分類部 204 は、それらを 1 つのグループに分類する。また、このようなインサイトサブジェクトの一部で、系列名が「売上」等の他の表記となっていた場合でも、表記統一部 203 が表記を統一するため、分類部 204 は、それらを 1 つのグループに分類することができる。

10

【0045】

ここで、上記のとおり、グループ化はインサイトタイプごとに行われる。よって、各インサイトタイプについて、グループ化の基準を予め定めておけばよい。インサイトタイプとしては、例えば相関が挙げられる。インサイトタイプが相関であるインサイトサブジェクトをグループ化する場合、分類部 204 は、相関関係の強さを評価できる、言い換えれば相関係数を計算可能なインサイトサブジェクトをグループ化すればよい。また、インサイトタイプが外れ値であるインサイトサブジェクトをグループ化する場合、分類部 204 は、外れ値を検出できるインサイトサブジェクト、つまり対応するデータ間の距離を計算可能なインサイトサブジェクトをグループ化すればよい。具体的には、例えば、分類部 204 は、各系列名を示す単語が同一のインサイトサブジェクトを 1 つのグループに分類してもよい。

20

【0046】

インサイトタイプとしては、相関以外にも任意のものを採用することができる。横断的複合インサイトを検出する場合、例えば、相互メジャー相関 (Cross-measure correlation)、二次元クラスタリング、帰属 (Attribution) 等のインサイトタイプを設定してもよい。

【0047】

また、例えば、分類部 204 は、シングルポイントインサイト (Single point insight)、すなわち 1 つのインサイトサブジェクトを入力とする横軸に順序が存在しない (non-ordinal dimension) インサイトサブジェクトをグループ化してもよい。このようなグループ化により、例えば、突出した No. 1 (Outstanding No.1)、突出した最下位 (Outstanding No. Last)、突出した上位 2 つ (Outstanding Top 2)、または均一度 (Evenness) 等のインサイトを検出することが可能になる。

30

【0048】

また、分類部 204 は、シングルシェープインサイト (Single shape insight)、すなわち 1 つのインサイトサブジェクトを入力とする横軸に順序が存在する (ordinal dimension) インサイトサブジェクトをグループ化してもよい。なお、横軸に順序が存在するデータとしては例えば時系列データが挙げられる。このようなグループ化により、変化点 (Change point)、トレンド、季節性 (Seasonality)、外れ値等のインサイトを検出することが可能になる。設定されるインサイトタイプには、横断的複合インサイトを検出可能なもの (例えば相関等) が少なくとも 1 つ含まれていればよく、横断的ではない複合インサイトを検出するためのもの (例えば変化点 (Change point) 等) が含まれていてもよい。

40

【0049】

粒度統一部 205 は、各インサイトサブジェクトにおけるデータの粒度を統一する。この処理は、評価部 206 がインサイトサブジェクト間の関連性を評価できるようにするための処理であるから、粒度が揃っていないデータを対象として行われる。粒度の統一は、

50

データセットから生成されたインサイトサブジェクトに対して行ってもよいし、分析対象となる複数のデータセットに対して予め行っておいてもよい。なお、データの粒度は、一連のデータがどのような細かさ（単位）であるかを示す。

【0050】

例えば、あるインサイトサブジェクトと他のインサイトサブジェクトが何れも月別の売上を示すものであるが、前者には毎月の売上が示されており、後者には隔月（奇数月）の売上が示されている場合、これらのデータの粒度は一致していない。この場合、両データ間の距離や類似度の評価ができないことがある。

【0051】

粒度統一部205は、このようなデータに対して粒度を揃える処理を行う。例えば、粒度統一部205は、欠損値補完によりデータを補完して粒度を揃えてもよいし、ダウンサンプリングにより粒度を揃えてもよい。欠損値補完は、他のデータから欠損部を予測して補完する処理であり、具体例としては内挿等が挙げられる。ダウンサンプリングは、サンプリング粒度を粗い方に合わせる処理である。

10

【0052】

上記の例において欠損値補完を行う場合、粒度統一部205は、他のインサイトサブジェクトにおける偶数月の売上を補完する。また、上記の例においてダウンサンプリングを行う場合、粒度統一部205は、あるインサイトサブジェクトにおける奇数月の売上のみが評価部206による評価に用いられるようにする。

【0053】

評価部206は、分類部204により同じグループに分類された複数のインサイトサブジェクトの組み合わせについてインサイトスコアを算出し、その算出結果を示す評価結果データ212を生成して記憶部21に記憶させる。例えば、評価部206は、同じグループに分類されたインサイトサブジェクトの組み合わせを入力としてインサイトスコアを返す関数 f_T を用いて上記の評価を行ってもよい。

20

【0054】

f_T は、インサイトタイプ T ごとに予め定義される関数であり、検出したいインサイトを与えるインサイトサブジェクトが入力されると高い値になるように設計される。インサイトタイプ T に対応するインサイトグループを G_T とすると、インサイトスコアは下記の式で表される。

30

【0055】

$$(\text{インサイトスコア}) = f_T(I_1, I_2, \dots, I_n | I_i \in G_T)$$

評価部206は、同じグループに分類された複数のインサイトサブジェクトを組にして、各組のインサイトスコアを算出してもよい。この場合、2つのインサイトサブジェクトを入力とする f_T を用いればよい。例えば、 $I_1 \sim I_3$ の3つのインサイトサブジェクトがグループ化されている場合、評価部206は、 I_1 と I_2 、 I_1 と I_3 、および I_2 と I_3 の各組をそれぞれ f_T に入力することにより、各組のインサイトスコアを算出する。

【0056】

インサイトスコアの算出方法は、インサイトタイプに応じたものとする。例えば、組にしたインサイトサブジェクト間の線形な相関の程度を評価する場合、評価部206は、ピアソン相関係数を算出する f_T を用いてインサイトスコアを算出してもよい。他にも、例えば、評価部206は、スピアマン順位相関係数やコサイン類似度、対応するデータ間のユークリッド距離やEMD (Earth Mover's distance) 等をインサイトスコアとして算出してもよい。

40

【0057】

なお、粒度統一部205がインサイトサブジェクトのデータの粒度を統一していた場合、評価部206は、粒度が統一された複数のインサイトサブジェクトの組み合わせについてインサイトスコアを算出する。異なるデータセット間では、データの粒度が不統一であることも多く、粒度が不統一であることが評価の支障となることも一般的には多いが、情報処理装置2によればそのような場合にも評価を行うことができる。すなわち、情報処理

50

装置 2 によれば、例示的实施形態 1 に係る情報処理装置 1 の奏する効果に加えて、粒度が不統一なデータを含むデータセットについても横断的複合インサイトを検出することが可能になるという効果が得られる。

【 0 0 5 8 】

出力データ生成部 2 0 7 は、評価結果データ 2 1 2 を用いて出力データ 2 1 3 を生成する。出力データ生成部 2 0 7 は、情報処理装置 2 の必須の構成要素ではないが、出力データ生成部 2 0 7 を設けることにより、情報処理装置 2 による分析の結果をより認識しやすい態様でユーザに提示することが可能になる。

【 0 0 5 9 】

(分析方法の流れ)

本例示的实施形態に係る分析方法の流れについて図 5 ~ 図 7 を参照して説明する。図 5 は、分析方法の流れを示すフロー図である。また、図 6 は、分析対象データ 2 1 1 と、当該分析対象データ 2 1 1 から生成されたインサイトサブジェクトの例を示す図である。そして、図 7 は、評価結果データ 2 1 2 と出力データ 2 1 3 の例を示す図である。

【 0 0 6 0 】

S 2 1 では、データ取得部 2 0 1 が、複数のデータセットの入力を受け付けて、分析対象データ 2 1 1 として記憶部 2 1 に記憶させる。例えば、データ取得部 2 0 1 は、入力部 2 3 を介して、図 6 に示す分析対象データ 2 1 1 の入力を受け付ける。分析対象データ 2 1 1 には、コンビニエンスストアにおける都道府県別の各月の売上を示すデータセット (D^S) と、スーパーマーケットにおける都道府県別の各月の売上を示すデータセット (D^T) が含まれる。

【 0 0 6 1 】

S 2 2 では、サブジェクト生成部 2 0 2 が、分析対象データ 2 1 1 に含まれる各データセットからインサイトサブジェクトを生成する。例えば、図 6 に示すデータセット D^S 、 D^T を用いる場合、サブジェクト生成部 2 0 2 は、データセット D^S からインサイトサブジェクト I^S_1 と I^S_2 を生成し、データセット D^T からインサイトサブジェクト I^T_1 と I^T_2 を生成することができる。

【 0 0 6 2 】

インサイトサブジェクト I^S_1 は、コンビニエンスストアにおける都道府県別の売上を示すものであり、図 6 では、 I^S_1 を売上の棒グラフ (横軸が都道府県、縦軸が売上) として示している。また、インサイトサブジェクト I^S_2 は、コンビニエンスストアにおける月毎の売上を示すものであり、図 6 では、 I^S_2 を売上の折れ線グラフ (横軸が日付、縦軸が売上) として示している。

【 0 0 6 3 】

同様に、インサイトサブジェクト I^T_1 は、スーパーマーケットにおける都道府県別の売上を示すものであり、図 6 では、 I^T_1 を売上の棒グラフ (横軸が都道府県、縦軸が売上) として示している。また、インサイトサブジェクト I^T_2 は、スーパーマーケットにおける月毎の売上を示すものであり、図 6 では、 I^T_2 を売上の折れ線グラフ (横軸が日付、縦軸が売上) として示している。

【 0 0 6 4 】

インサイトサブジェクト I は、例えば下記のようなデータ形式とすることもできる。

$I = \{ \text{subspace, breakdown, measure, aggregation} \}$

上記 “subspace” (サブスペース) は、多次元データであるデータセットに含まれるレコードをどのようにフィルタしたかを示す。上記 “subspace” は、各チャートの凡例に対応する。例えば、図 6 の I^S_2 の折れ線グラフにおける “subspace” は「東京都」である。フィルタリングを行わないことは、“*” 等の記号で表せばよい。

【 0 0 6 5 】

上記 “breakdown” (ブレイクダウン) は、多次元データであるデータセットを集計するキーとして使用されるカラムを示す。上記 “breakdown” は、各チャートの横軸に対応する。例えば、図 6 の I^S_2 の折れ線グラフにおける “breakdown” は「日付」である。

10

20

30

40

50

【 0 0 6 6 】

上記 “ measure ” (メジャー) は、多次元データであるデータセットにおいて数値データとして使用されるカラムを示す。上記 “ measure ” は、各チャートの縦軸に対応する。例えば、図 6 の I^S_2 の折れ線グラフにおける “ measure ” は「売上」の数値データである。

【 0 0 6 7 】

上記 “ aggregation ” (アグリゲーション) は、“ breakdown ” ごとにデータを集計する方法 (例えば関数) を示す。上記 “ aggregation ” の例としては、合計、平均、最大値、最小値等が挙げられる。集計に用いられる関数が「合計」である場合、“ aggregation ” は省略してもよい。

【 0 0 6 8 】

例えば、図 6 に示す I^S_2 であれば、 $I^S_2 = \{ \{ *, 東京都 \}, 日付, 売上 \}$ と表すことができる。 S_{22} では、サブジェクト生成部 202 は、分析対象データ 211 に含まれる各データセットからこのようなデータ形式のインサイトサブジェクトを生成してもよい。

10

【 0 0 6 9 】

S_{23} では、表記統一部 203 が、 S_{22} で生成された各インサイトサブジェクトにおけるデータの表記を統一する。例えば、図 6 に示す I^S_1 、 I^S_2 、 I^T_1 、 I^T_2 の中では、 I^S_1 における横軸のラベル「都道府県」と、 I^T_1 における横軸のラベル「場所」の意味が類似している。また、 I^S_1 の系列名「東京都」、「大阪府」、「神奈川県」は、 I^T_1 の系列名「東京」、「大阪」、「神奈川」のそれぞれと意味および表記が類似している。表記統一部 203 は、このような単語を抽出し、それらの表記を統一する。例えば、表記統一部 203 は、 I^S_1 における横軸のラベルを「場所」に置換し、系列名「東京都」、「大阪府」、「神奈川県」を、それぞれ「東京」、「大阪」、「神奈川」に置換してもよい。

20

【 0 0 7 0 】

S_{24} では、分類部 204 が、 S_{22} で生成されたインサイトサブジェクトであって、 S_{23} で表記が統一されたインサイトサブジェクトをグループ化する。例えば、図 6 に示す I^S_1 、 I^S_2 、 I^T_1 、 I^T_2 のうち、縦軸と横軸のラベルが共通するインサイトサブジェクトをグループ化するとする。この場合、分類部 204 は、縦軸のラベルが「売上」で横軸のラベルが「場所」である I^S_1 と I^T_1 をグループ化する。 I^S_1 の「都道府県」は表記統一部 203 により「場所」に置換済みであるからこのようなグループ化が可能になっている。また、分類部 204 は、縦軸のラベルが「売上」で横軸のラベルが「日付」である I^S_2 と I^T_2 をグループ化する。

30

【 0 0 7 1 】

I^S_1 と I^T_1 を含むグループを G^1 、 I^S_2 と I^T_2 を含むグループを G^2 とすると、グループ化の結果は下記のように表される。

$$\begin{matrix} I^S_1, I^T_1 & G^1 \\ I^S_2, I^T_2 & G^2 \end{matrix}$$

S_{25} では、粒度統一部 205 が、 S_{24} でグループ化されたインサイトサブジェクトに含まれるデータの粒度を統一する。例えば、図 6 に示す I^S_2 の「日付」は、奇数月の 1 日であるのに対し、 I^T_2 の「日付」は毎月の 1 日である。粒度統一部 205 は、このように粒度に差異があるデータを抽出し、それらのデータの粒度を揃える処理を行う。例えば、粒度統一部 205 は、 I^T_2 の「日付」のデータのうち、奇数月のデータを抽出 (すなわちダウンサンプリング) することにより、「日付」データの粒度を揃えてもよい。また、粒度統一部 205 は、 I^S_2 の偶数月のデータを欠損値補完することにより、「日付」データの粒度を揃えてもよい。なお、欠損値補完は、データのサンプリング日付にずれがある場合にも有効である。例えば、粒度統一部 205 は、毎月 1 日のデータと、毎月 15 日のデータの粒度を揃える場合、毎月 15 日のデータを欠損値補完することにより、毎月 1 日のデータを生成してもよい。

40

【 0 0 7 2 】

S_{26} では、評価部 206 が、 S_{24} でグループ化され、 S_{25} でデータの粒度が統一

50

されたインサイトサブジェクトの組み合わせを評価し、評価結果を評価結果データ 2 1 2 として記憶部 2 1 に記憶させる。より詳細には、評価部 2 0 6 は、同じグループに含まれるインサイトサブジェクトを組にして、その組についてのインサイトスコアを算出する、という処理を各グループについて行う。

【 0 0 7 3 】

例えば、評価部 2 0 6 は、 $f_T(I_i, I_j)$ の式で表されるスコア関数、すなわち評価対象とする 2 つのインサイトサブジェクトを入力とし、インサイトスコアを出力とする関数を用いてインサイトスコアを算出してもよい。このスコア関数を用いる場合、グループ G^1 のインサイトスコアは $f_T(I^{S_1}, I^{T_1})$ 、グループ G^2 のインサイトスコアは $f_T(I^{S_2}, I^{T_2})$ と表される。

10

【 0 0 7 4 】

評価部 2 0 6 は、上述のような評価結果をリスト化することにより、例えば図 7 に示すような評価結果データ 2 1 2 を生成してもよい。図 7 に示す評価結果データ 2 1 2 は、インサイトサブジェクトの組み合わせと、その組み合わせについて算出されたインサイトスコアとを示すテーブル形式のデータである。また、図 7 に示す評価結果データ 2 1 2 には、インサイトスコアの順位を示す「ランク」と、「インサイトタイプ」についても示されている。このように、評価部 2 0 6 は、インサイトサブジェクトの組み合わせと、その組み合わせについて算出されたインサイトスコアに加えて、評価に関する各種情報を含む評価結果データ 2 1 2 を生成してもよい。

【 0 0 7 5 】

S 2 7 では、出力データ生成部 2 0 7 が、S 2 6 で生成された評価結果データ 2 1 2 を用いて出力データ 2 1 3 を生成し、出力部 2 4 に出力させる。例えば、図 7 に示す評価結果データ 2 1 2 を用いる場合、出力データ生成部 2 0 7 は、インサイトスコア（ランク）が最も高いインサイトサブジェクトの組み合わせを示す出力データ 2 1 3 を生成し、出力部 2 4 に出力させる。これにより、図 5 の処理は終了する。

20

【 0 0 7 6 】

出力データ 2 1 3 は、インサイトをユーザが認識しやすいように、当該インサイトを可視化したものであってもよい。可視化方法は、インサイトタイプに応じて決定すればよい。例えば、出力データ生成部 2 0 7 は、インサイトタイプが「相関」である場合、インサイトに関する情報として相関関係を表すのに適したチャート（例えば二次元の散布図）を出力データ 2 1 3 として生成してもよい。

30

【 0 0 7 7 】

図 7 の下側には、評価結果データ 2 1 2 に示されるインサイトサブジェクトの組み合わせのうち、最もインサイトスコアが高かった（つまり、ランクが 1 の）ものについてのインサイトに関する情報の例を示している。具体的には、図 7 に示されるインサイトに関する情報には、スーパーマーケットとコンビニエンスストアの売上の相関を示す散布図と、インサイトの詳細を示すインサイト情報とが含まれている。インサイト情報には、インサイトタイプとインサイトスコアの他、各インサイトサブジェクトの詳細とその元になったデータセットが示されている。このような情報を出力部 2 4 に出力させることにより、情報処理装置 2 のユーザに、スーパーマーケットとコンビニエンスストアの売上の推移に強い相関がある、というインサイトを容易に認識させることができる。

40

【 0 0 7 8 】

無論、出力データ生成部 2 0 7 が生成する情報は、インサイトをユーザに認識させることができるようなものであればよく、図 7 の例に限られない。例えば、出力データ生成部 2 0 7 は、最もインサイトスコアが高かったインサイトサブジェクトの組み合わせについて、各インサイトサブジェクトのチャートを生成し、これを出力データ 2 1 3 としてもよい。

【 0 0 7 9 】

なお、分析結果をユーザに提示する際に、必ずしも新たな出力データ 2 1 3 を生成する必要はない。例えば、評価部 2 0 6 が、図 7 に示す評価結果データ 2 1 2 の全部または一

50

部を出力部 2 4 に出力させることにより、分析結果をユーザに提示してもよい。また、評価部 2 0 6 は、ランクが 1 となった各インサイトサブジェクトや、インサイトスコアが所定の閾値以上となった各インサイトサブジェクトを構成するデータを出力させてもよい。このように、分析結果を出力させる態様は任意であり、図 7 のような例に限定されない。また、分析結果の可視化方法をユーザに選択させてもよい。この場合、出力データ生成部 2 0 7 は、ユーザが選択した方法で分析結果を可視化する。

【 0 0 8 0 】

このように、情報処理装置 2 は、複数のデータセットの分析結果として、インサイトの発見に繋がる可能性のあるチャートやデータ等を出力することができる。これにより、人手でチャートを比較する必要がなくなる。また、最終的にはインサイトをユーザが検討する場合であっても、分析に役立つようなデータセットを容易に絞り込むことができる。よって、分析・可視化に要する時間を大幅に短縮することができる。

10

【 0 0 8 1 】

また、情報処理装置 2 を用いることにより、全ての分析をユーザが行う場合に生じる判断基準のブレが発生する余地もない。さらに、分析をユーザが行う場合に生じる見逃しのリスク等も低減することができる。また、大規模なデータセットが分析対象である場合、ユーザによる複合インサイトの発見は困難であるが、情報処理装置 2 によれば、複合インサイト（横断的複合インサイトも含む）の発見が容易になる。

【 0 0 8 2 】

なお、図 5 のフローチャートにおいて、S 2 3 の処理は、S 2 4 の処理よりも先に行えばよく、例えば S 2 1 と S 2 2 の間に行ってもよい。また、S 2 5 の処理は、S 2 6 の処理よりも先に行えばよく、例えば S 2 1 と S 2 2 の間に行ってもよい。

20

【 0 0 8 3 】

（粒度の違いへの対応の変形例）

評価部 2 0 6 は、データの粒度が異なる複数のインサイトサブジェクトの組み合わせについてもインサイトスコアを算出可能な評価方法により、インサイトサブジェクトを評価してもよい。これにより、例示的实施形態 1 に係る情報処理装置 1 の奏する効果に加えて、粒度が不統一なデータを含むデータセットについても横断的複合インサイトを検出することが可能になるという効果が得られる。また、この場合、粒度統一部 2 0 5 を省略することができるという効果も得られる。

30

【 0 0 8 4 】

例えば、インサイトサブジェクトにおける横軸のデータに順序が存在する（ordinal dimension である）場合には、評価部 2 0 6 は、DTW（Dynamic Time Warping：動的時間伸縮法）や関数データ解析によりインサイトスコアを算出してもよい。なお、順序が存在するデータの例としては、例えば時系列データ等が挙げられる。DTWでは、 $s = (s_1, \dots, s_n)$ と $t = (t_1, \dots, t_m)$ の要素間の距離を総当りで計算したコスト行列 W の端 $(1, 1)$ から端 (n, n) の最短経路を動的計画法で求める。DTWによれば、サンプルサイズが異なるデータ間の距離や類似度を計算可能であり、そのような距離や類似度をインサイトスコアの計算に用いることができる。また、関数データ解析を用いる場合、評価部 2 0 6 は、各インサイトサブジェクトのレコードを表現する連続的な関数を導出し、その関数を介してインサイトサブジェクト間の距離や類似度を計算し、それらをインサイトスコアの計算に用いることができる。

40

【 0 0 8 5 】

〔例示的实施形態 3〕

本発明の第 3 の例示的实施形態について、図面を参照して詳細に説明する。上述の例示的实施形態において、インサイトサブジェクトをグループ化したときに、3 つ以上のインサイトサブジェクトが 1 つのグループに分類されることがあり得る。このような場合、上述したスコア関数 $f_T(I_i, I_j)$ では、3 つ以上のインサイトサブジェクトをまとめて評価することはできない。また、3 つ以上のインサイトサブジェクトをまとめて評価する方法については、特許文献 1 にも記載も示唆もされていない。

50

【 0 0 8 6 】

本例示的实施形態では、3つ以上のインサイトサブジェクトをまとめて評価することが可能な評価方法について図8～図10に基づいて説明する。図8は、本例示的实施形態に係る情報処理装置3の構成を示すブロック図である。図9は、本例示的实施形態に係る分析方法の流れを示すフロー図である。図10は、インサイトスコアの算出方法と、外れ値の検出方法を説明する図である。

【 0 0 8 7 】

(情報処理装置3の構成)

図8に示すように、情報処理装置3は、評価部31と外れ値検出部32を備えている。なお、外れ値を検出する必要がない場合には外れ値検出部32を省略してもよい。評価部31は、図1に示した評価部12および図4に示した評価部206と同様に、グループ化された複数のインサイトサブジェクトの組み合わせについてインサイトスコアを算出する。評価部31は、3つ以上のインサイトサブジェクトをまとめて評価することができる点、言い換えれば3つ以上のインサイトサブジェクトにおけるインサイトの有無を示す1つのインサイトスコアを算出できる点で、評価部12、206と相違している。

10

【 0 0 8 8 】

具体的には、評価部31は、グループ化された複数のインサイトサブジェクトを主成分分析することにより求めた、各主成分の寄与度の偏りの程度に基づいて当該インサイトサブジェクトの組み合わせについてのインサイトスコアを算出する。主成分分析は、任意の数のインサイトサブジェクトを対象として行うことができる。このため、本例示的实施形態に係る情報処理装置3によれば、例示的实施形態1、2に係る情報処理装置1、2の奏する効果に加えて、3つ以上のインサイトサブジェクトをまとめて評価することが可能になるという効果が得られる。なお、評価方法の詳細およびこのような評価が可能である理由については、図9および図10に基づいて後述する。

20

【 0 0 8 9 】

外れ値検出部32は、評価部31による主成分分析により求められた主成分を用いて、グループ化された複数のインサイトサブジェクトに含まれるデータを表すことにより、当該データに含まれる外れ値を検出する。このため、本例示的实施形態に係る情報処理装置3によれば、例示的实施形態1、2に係る情報処理装置1、2の奏する効果に加えて、評価のために行った主成分分析の結果を利用した効率のよい外れ値検出ができるという効果が得られる。なお、外れ値検出方法の詳細およびこのような方法で外れ値を検出することが可能である理由については、図9および図10に基づいて後述する。

30

【 0 0 9 0 】

(情報処理装置3が実行する処理の流れ)

情報処理装置3が実行する処理の流れを図9に基づいて説明する。なお、図9の処理の前に、複数のインサイトサブジェクトがグループ化済であるとする。つまり、図8には示していないが、本例示的实施形態では、情報処理装置3が分類部11(例示的实施形態1)または分類部204(例示的实施形態2)に相当する構成を備えていることを想定している。なお、情報処理装置3は、情報処理装置2が備える各種構成(例えば、データ取得部201やサブジェクト生成部202等)の一部または全部を備えていてもよい。

40

【 0 0 9 1 】

S31では、評価部31が、インサイトサブジェクトのグループを評価する。より詳細には、まず、評価部31は、評価対象のグループに含まれる各インサイトサブジェクトにおける、主成分分析の対象とするデータを特定する。例えば、インサイトサブジェクトが $I = \{ \text{subspace, breakdown, measure, aggregation} \}$ の形式で表されていた場合、評価部31は、各インサイトサブジェクトにおける“measure”の項目のデータを主成分分析の対象とすればよい。

【 0 0 9 2 】

次に、評価部31は、主成分分析の対象として特定したデータについて主成分分析を行う。例えば、評価部31は、各インサイトサブジェクトにおける“measure”の項目のデー

50

タから多次元の相関行列を生成し、この相関行列を用いて主成分分析を行ってもよい。主成分分析により、固有値と固有ベクトルが算出される。

【 0 0 9 3 】

続いて、評価部 3 1 は、算出された固有値を用いて、各主成分の寄与率を算出する。各主成分の寄与率はその軸方向（固有ベクトル）における情報量とみなすことができるから、各主成分の寄与率の偏り度合いを調べることで、インサイトサブジェクト間の相関の強さを定量的に評価することができる。

【 0 0 9 4 】

例えば、図 1 0 には、相関がないインサイトサブジェクトを主成分分析して算出された各主成分の寄与率を示す棒グラフ 1 0 0 1 と、相関があるインサイトサブジェクトを主成分分析して算出された各主成分の寄与率を示す棒グラフ 1 0 0 2 を示している。なお、図 1 0 において、P C 1 は第 1 主成分、P C 2 は第 2 主成分、P C 3 は第 3 主成分である。

【 0 0 9 5 】

棒グラフ 1 0 0 1 では、P C 1 ~ P C 3 の寄与率は概ね同程度であり、主成分間での偏り度合いは小さい。一方、棒グラフ 1 0 0 2 では、P C 1 の寄与率が最も高く、P C 2 の寄与率はその半分程度であり、P C 3 の寄与率はかなり小さく、全体として偏り度合いが大きい。このように、インサイトサブジェクト間の相関の有無は、各主成分の寄与率の偏り度合いに明瞭に反映される。

【 0 0 9 6 】

したがって、各主成分の寄与率の偏り度合いを定量的に評価すれば、その評価結果をインサイトスコアとすることができる。例えば、第 1 主成分の寄与率をインサイトスコアとしてもよい。これは、図 1 0 に示されるように、各主成分の寄与率の偏り度合いが大きい場合（棒グラフ 1 0 0 2）には、小さい場合（棒グラフ 1 0 0 1）と比べて第 1 主成分 P C 1 の寄与率が大きいためである。

【 0 0 9 7 】

また、図 1 0 に示されるように、各主成分の寄与率の偏り度合いが大きい場合（棒グラフ 1 0 0 2）には、P C 1 ~ P C 3 の中で寄与率が突出して高いもの（具体的には P C 1）が存在する。一方、各主成分の寄与率の偏り度合いが小さい場合（棒グラフ 1 0 0 1）には、寄与率が突出して高いものは存在しない。このため、例えば、各主成分の寄与率を入力とし、入力された寄与率の中に突出して高いものが含まれているほど高い値を出力するスコア関数を用いてインサイトスコアを算出することもできる。

【 0 0 9 8 】

なお、インサイトサブジェクト間の非線形な相関を検出したい場合には、評価部 3 1 は、通常の主成分分析のかわりに、任意のカーネルを用いたカーネル主成分分析を実行してもよい。また、レコードのサンプリング粒度の違いなどで相関行列が計算できない場合には、評価部 3 1 は、関数データ解析を用いた関数主成分分析を実行してもよい。

【 0 0 9 9 】

S 3 2 では、外れ値検出部 3 2 が、グループ化された各インサイトサブジェクトに含まれる外れ値の検出を行う。例えば、S 3 1 で各インサイトサブジェクトにおける“measure”の項目のデータを用いた評価が行われていた場合、外れ値検出部 3 2 も各インサイトサブジェクトにおける“measure”の項目のデータにおける外れ値を検出する。

【 0 1 0 0 】

外れ値の検出は、S 3 1 における評価のために行われた主成分分析により求められた主成分を用いて、グループ化された複数のインサイトサブジェクトに含まれるデータを表すことにより行われる。

【 0 1 0 1 】

図 1 0 の 1 0 0 3 は、サンプルデータを主成分分析して求めた第 1 主成分 P C 1 と第 2 主成分 P C 2 により当該サンプルデータを表した点を、縦軸を P C 2、横軸を P C 1 とする座標平面上にプロットしたものである。主成分分析後のプロットにおいて、他のデータと離れているデータは、元のサンプルデータにおいても他のデータと離れている。よって

10

20

30

40

50

、1003において「外れ値」とされているプロットのように、他のデータから離れたデータを外れ値として検出すればよい。

【0102】

例えば、外れ値検出部32は、主成分で表されたデータのHotellingの T^2 統計量を算出し、算出した T^2 統計量が顕著なデータを外れ値として検出してもよい。図10の1004は、同図の1003に示すサンプルデータから算出した T^2 統計量を、横軸がサンプル番号、縦軸が T^2 統計量の座標平面にプロットしたものである。同図の1003において「外れ値」とされていたプロットは、 T^2 統計量が他のプロットと比べて大きい値となっている。よって、外れ値検出部32は、 T^2 統計量を用いて外れ値を検出することができる。

10

【0103】

また、 T^2 統計量はF分布や χ^2 分布に従うことが知られている。このため、外れ値検出部32は、統計的検定に基づいて得られたp値を用いてスコアを計算してもよい。この場合、外れ値検出部32は、算出したスコアを用いて外れ値を検出すればよい。

【0104】

以上により、図9の処理は終了する。なお、S31の評価結果とS32で検出された外れ値は、評価結果データとして記憶しておけばよい。評価結果データは、そのまま出力してもよいし、例示的实施形態2と同様に、評価結果データから出力データを生成し、生成した出力データを出力してもよい。

【0105】

〔参考例〕

評価部31による上述の評価方法は、横断的複合インサイトの検出に好適であると共に、横断的ではない、つまり1つのデータセットにおけるインサイトの検出にも好適である。このため、上述の情報処理装置3は、必ずしも分類部204（例示的实施形態2）や、分類部11（例示的实施形態1）に相当する構成を備えている必要はない。

20

【0106】

本参考例に係る情報処理装置3は、評価対象となる複数のインサイトサブジェクトを取得する取得部と、上述の評価部31を備えている。前記取得部が取得する複数のインサイトサブジェクトは、少なくとも1つのデータセットから生成されたものであればよい。つまり、複数のデータセットから生成された複数のインサイトサブジェクトを用いることが必須ではない点で、本参考例と上述の各例示的实施形態は相違している。

30

【0107】

本参考例の情報処理装置によれば、評価部31は、取得部が取得した複数の前記インサイトサブジェクトを主成分分析することにより得られた、各主成分の寄与度の偏りの程度に基づいて、当該インサイトサブジェクトの組み合わせについてのインサイトスコアを算出する。よって、3つ以上のインサイトサブジェクトをまとめて評価することができなかったという従来の課題を解決することができる。

【0108】

また、本参考例に係る分析方法は、少なくとも1つのプロセッサが、評価対象となる複数のインサイトサブジェクトを取得すること、および、取得した複数の前記インサイトサブジェクトを主成分分析することにより得られた、各主成分の寄与度の偏りの程度に基づいて、当該インサイトサブジェクトの組み合わせについてのインサイトスコアを算出すること、を含む。そして、本参考例に係る分析プログラムは、コンピュータに、評価対象となる複数のインサイトサブジェクトを取得する処理と、取得した複数の前記インサイトサブジェクトを主成分分析することにより得られた、各主成分の寄与度の偏りの程度に基づいて、当該インサイトサブジェクトの組み合わせについてのインサイトスコアを算出する処理と、を実行させる。これらの分析方法および分析プログラムによっても、3つ以上のインサイトサブジェクトをまとめて評価することができなかったという従来の課題を解決することができる。

40

【0109】

50

〔変形例〕

上述の例示的实施形態 1 において、1つの情報処理装置 1 が行っていた処理は、複数の情報処理装置に分担させてもよい。言い換えれば、情報処理装置 1 が行う処理の一部を、少なくとも1つの他の情報処理装置に実行させてもよい。さらに言い換えれば、上述の各処理を少なくとも1つのプロセッサに行わせる場合、その少なくとも1つのプロセッサは、1つの情報処理装置 1 が備えているものであってもよいし、それぞれ異なる情報処理装置が備えているものであってもよい。これは、上述の例示的实施形態 2 における情報処理装置 2、および例示的实施形態 3 における情報処理装置 3 についても同様である。

【0110】

〔ソフトウェアによる実現例〕

情報処理装置 1 ~ 3 の一部又は全部の機能は、集積回路 (IC チップ) 等のハードウェアによって実現してもよいし、ソフトウェアによって実現してもよい。

【0111】

後者の場合、情報処理装置 1 ~ 3 は、例えば、各機能を実現するソフトウェアであるプログラムの命令を実行するコンピュータによって実現される。このようなコンピュータの一例 (以下、コンピュータ C と記載する) を図 11 に示す。コンピュータ C は、少なくとも1つのプロセッサ C 1 と、少なくとも1つのメモリ C 2 と、を備えている。メモリ C 2 には、コンピュータ C を情報処理装置 1 ~ 3 として動作させるためのプログラム P が記録されている。コンピュータ C において、プロセッサ C 1 は、プログラム P をメモリ C 2 から読み取って実行することにより、情報処理装置 1 ~ 3 の各機能が実現される。

【0112】

プロセッサ C 1 としては、例えば、CPU (Central Processing Unit)、GPU (Graphic Processing Unit)、DSP (Digital Signal Processor)、MPU (Micro Processing Unit)、FPU (Floating point number Processing Unit)、PPU (Physics Processing Unit)、マイクロコントローラ、又は、これらの組み合わせなどを用いることができる。メモリ C 2 としては、例えば、フラッシュメモリ、HDD (Hard Disk Drive)、SSD (Solid State Drive)、又は、これらの組み合わせなどを用いることができる。

【0113】

なお、コンピュータ C は、プログラム P を実行時に展開したり、各種データを一時的に記憶したりするための RAM (Random Access Memory) を更に備えていてもよい。また、コンピュータ C は、他の装置との間でデータを送受信するための通信インタフェースを更に備えていてもよい。また、コンピュータ C は、キーボードやマウス、ディスプレイやプリンタなどの入出力機器を接続するための入出力インタフェースを更に備えていてもよい。

【0114】

また、プログラム P は、コンピュータ C が読み取り可能な、一時的でない有形の記録媒体 M に記録することができる。このような記録媒体 M としては、例えば、テープ、ディスク、カード、半導体メモリ、又はプログラマブルな論理回路などを用いることができる。コンピュータ C は、このような記録媒体 M を介してプログラム P を取得することができる。また、プログラム P は、伝送媒体を介して伝送することができる。このような伝送媒体としては、例えば、通信ネットワーク、又は放送波などを用いることができる。コンピュータ C は、このような伝送媒体を介してプログラム P を取得することもできる。

【0115】

〔付記事項 1〕

本発明は、上述した実施形態に限定されるものでなく、請求項に示した範囲で種々の変更が可能である。例えば、上述した実施形態に開示された技術的手段を適宜組み合わせて得られる実施形態についても、本発明の技術的範囲に含まれる。

【0116】

〔付記事項 2〕

10

20

30

40

50

上述した実施形態の一部又は全部は、以下のようにも記載され得る。ただし、本発明は、以下の記載する態様に限定されるものではない。

【0117】

(付記1)

複数のデータセットのそれぞれから当該データセットに含まれる複数のデータ項目を関連付けることにより生成されたデータであるインサイトサブジェクトを、検出対象のインサイトごとにグループ化する分類手段と、グループ化された複数の前記インサイトサブジェクトの組み合わせについて、インサイトの有無を判定するための評価値を算出する評価手段と、を備える情報処理装置。この構成によれば、複数のデータセット間におけるインサイトの検出を可能にすることができる。

10

【0118】

(付記2)

複数の前記インサイトサブジェクトにおける表記を統一する表記統一手段をさらに備え、前記分類手段は、表記が統一された前記インサイトサブジェクトをグループ化する、付記1に記載の情報処理装置。この構成によれば、表記が不統一なデータセットについても横断的複合インサイトを検出することが可能になる。

【0119】

(付記3)

複数の前記インサイトサブジェクトにおけるデータの粒度を統一する粒度統一手段をさらに備え、前記評価手段は、粒度が統一された複数の前記インサイトサブジェクトについて前記評価値を算出する、付記1または2に記載の情報処理装置。この構成によれば、粒度が不統一なデータを含むデータセットについても横断的複合インサイトを検出することが可能になる。

20

【0120】

(付記4)

前記評価手段は、動的時間伸縮法または関数データ解析により前記評価値を算出する、付記1または2に記載の情報処理装置。この構成によれば、粒度が不統一なデータを含むデータセットについても横断的複合インサイトを検出することが可能になる。

【0121】

(付記5)

前記評価手段は、グループ化された複数の前記インサイトサブジェクトを主成分分析することにより求めた、各主成分の寄与度の偏りの程度に基づいて前記評価値を算出する、付記1から4の何れかに記載の情報処理装置。この構成によれば、3つ以上のインサイトサブジェクトをまとめて評価することが可能になる。

30

【0122】

(付記6)

前記主成分分析により求められた主成分を用いて、グループ化された複数の前記インサイトサブジェクトに含まれるデータを表すことにより、当該データに含まれる外れ値を検出する外れ値検出手段をさらに備える、付記5に記載の情報処理装置。この構成によれば、評価のために行った主成分分析の結果を利用した効率のよい外れ値検出ができる。

40

【0123】

(付記7)

少なくとも1つのプロセッサが、複数のデータセットのそれぞれから当該データセットに含まれる複数のデータ項目を関連付けることにより生成されたデータであるインサイトサブジェクトを、検出対象のインサイトごとにグループ化すること、およびグループ化された複数の前記インサイトサブジェクトの組み合わせについて、インサイトの有無を判定するための評価値を算出すること、を含む分析方法。この構成によれば、複数のデータセット間におけるインサイトの検出を可能にすることができる。

【0124】

(付記8)

50

コンピュータに、複数のデータセットのそれぞれから当該データセットに含まれる複数のデータ項目を関連付けることにより生成されたデータであるインサイトサブジェクトを、検出対象のインサイトごとにグループ化する処理と、グループ化された複数の前記インサイトサブジェクトの組み合わせについて、インサイトの有無を判定するための評価値を算出する処理と、を実行させる分析プログラム。この構成によれば、複数のデータセット間におけるインサイトの検出を可能にすることができる。

【0125】

(付記9)

少なくとも1つのプロセッサを備え、前記プロセッサは、複数のデータセットのそれぞれから当該データセットに含まれる複数のデータ項目を関連付けることにより生成されたデータであるインサイトサブジェクトを、検出対象のインサイトごとにグループ化する処理と、グループ化された複数の前記インサイトサブジェクトの組み合わせについて、インサイトの有無を判定するための評価値を算出する処理とを実行する情報処理装置。

10

【0126】

なお、この情報処理装置は、更にメモリを備えていてもよく、このメモリには、前記をグループ化する処理と、前記評価する処理とを前記プロセッサに実行させるためのプログラムが記憶されていてもよい。また、このプログラムは、コンピュータ読み取り可能な一時的でない有形の記録媒体に記録されていてもよい。

【符号の説明】

【0127】

20

- 1 情報処理装置
- 1 1 分類部(分類手段)
- 1 2 評価部(評価手段)
- 2 情報処理装置
- 2 0 3 表記統一部(表記統一手段)
- 2 0 4 分類部(分類手段)
- 2 0 5 粒度統一部(粒度統一手段)
- 2 0 6 評価部(評価手段)
- 3 情報処理装置
- 3 1 評価部(評価手段)
- 3 2 外れ値検出部(外れ値検出手段)

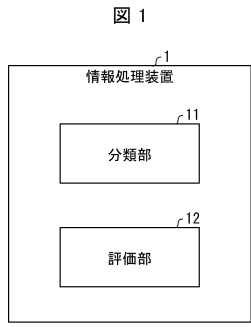
30

40

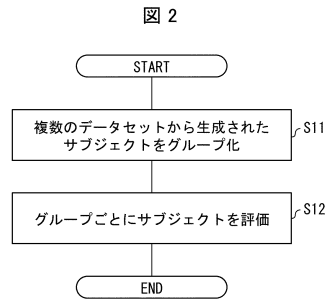
50

【図面】

【図 1】

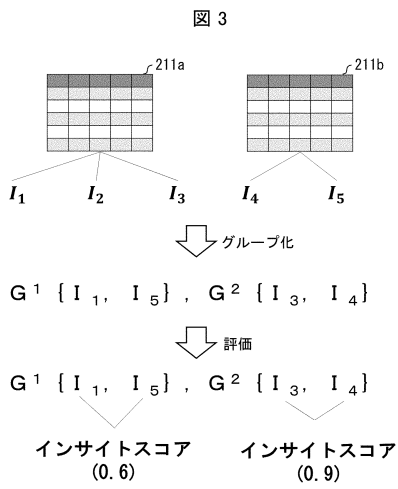


【図 2】

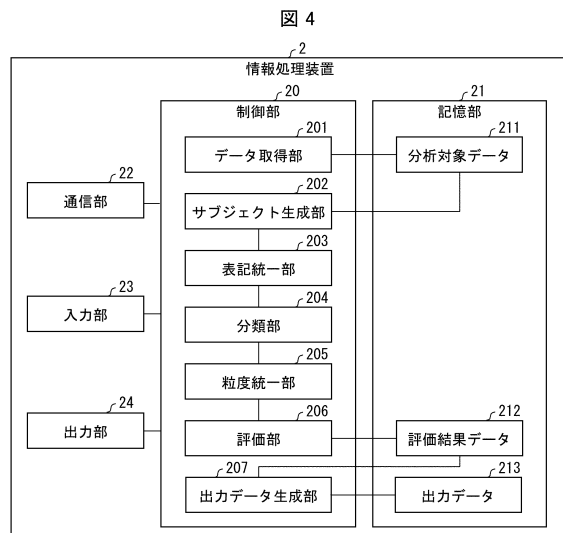


10

【図 3】



【図 4】



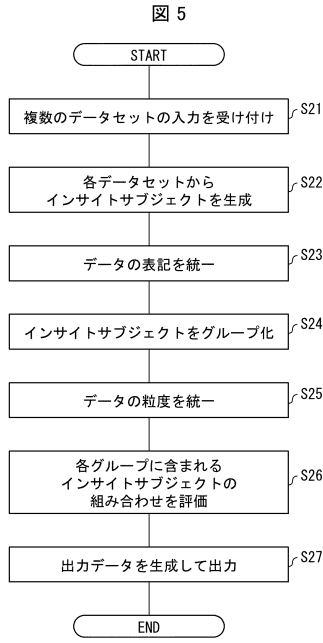
20

30

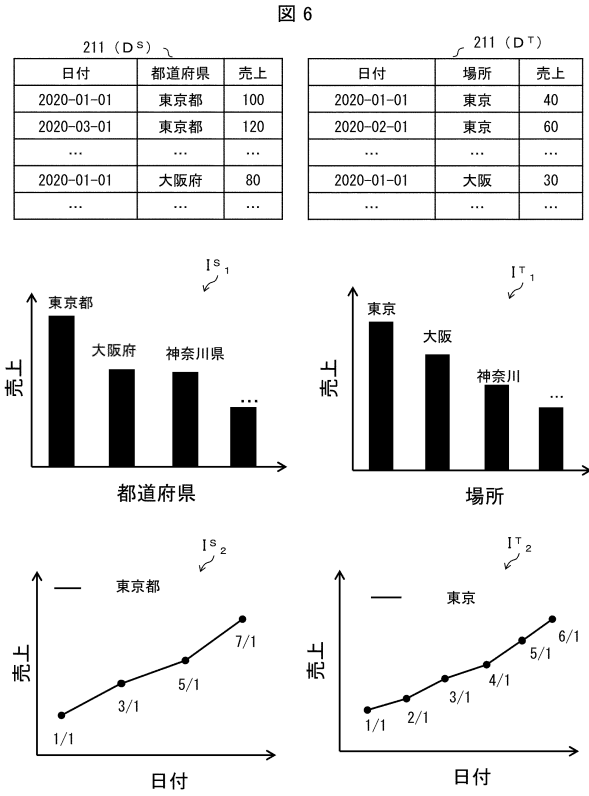
40

50

【 図 5 】



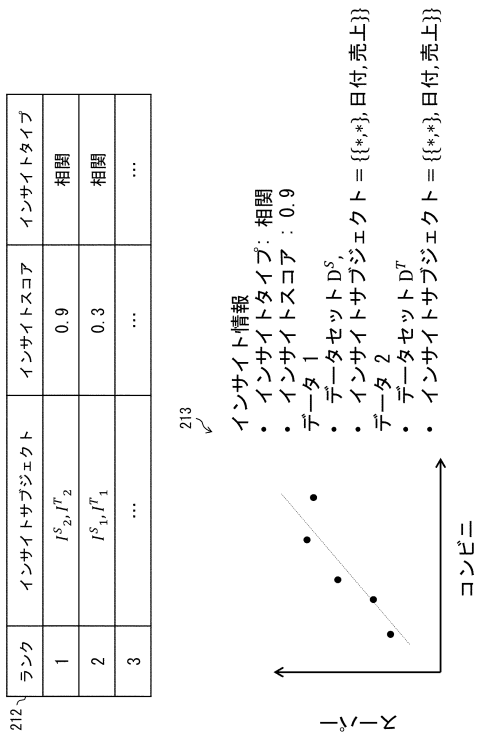
【 図 6 】



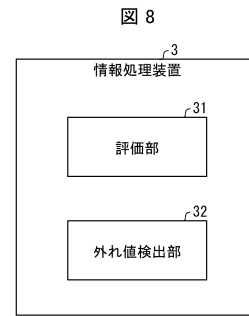
10

20

【 図 7 】



【 図 8 】

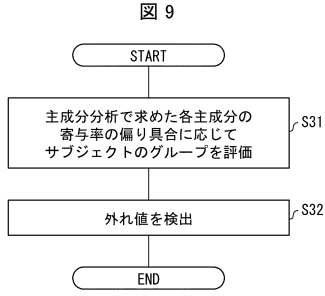


30

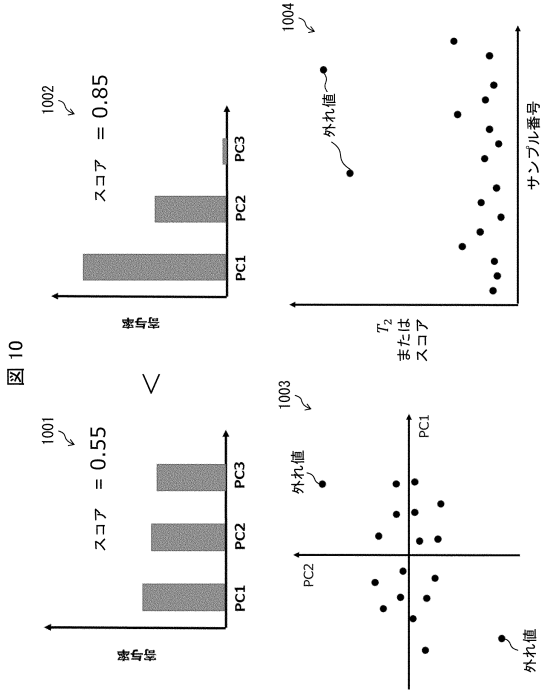
40

50

【 図 9 】



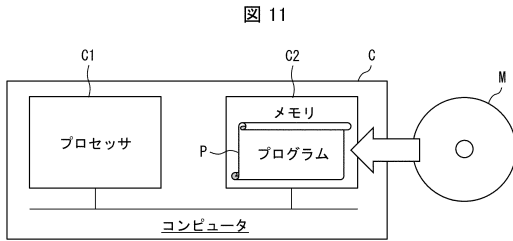
【 図 10 】



10

20

【 図 11 】



30

40

50

フロントページの続き

- 株式会社内
(72)発明者 草野 元紀
東京都港区芝五丁目7番1号 日本電気株式会社内
- 審査官 山崎 誠也
- (56)参考文献 国際公開第2022/026378(WO, A1)
特開2020-187511(JP, A)
米国特許出願公開第2020/0257682(US, A1)
特開2019-148897(JP, A)
特開2021-043899(JP, A)
国際公開第2017/163277(WO, A1)
塚越 雄登, 次元間の関係に着目したドメインオンロジーに基づく異種データ間の関連性
発見, 情報処理学会 研究報告 知能システム(ICS), 日本, 情報処理学会, 2020年09
月07日, p.1-8, ISSN:2188-885X
NEC、AIで予測分析した結果を可視化し、次の一手を提示するdotDataの新サービスを販売
開始, [online], 2020年10月07日, p.1-3, [2022年1月6日検索], インターネット URL:htt
ps://web.archive.org/web/20201124042138/https://jpn.nec.com/press/202010/202010
07_01.html
- (58)調査した分野 (Int.Cl., DB名)
G06Q 10/00-99/00