

(19) 日本国特許庁(JP)

(12) 公開特許公報(A)

(11) 特許出願公開番号

特開2019-56983

(P2019-56983A)

(43) 公開日 平成31年4月11日(2019.4.11)

(51) Int.Cl.	F I	テーマコード (参考)
GO6F 16/00 (2019.01)	GO6F 17/30	220Z
GO6N 20/00 (2019.01)	GO6F 17/30	210D
	GO6N 99/00	150

審査請求 未請求 請求項の数 9 O L (全 31 頁)

(21) 出願番号	特願2017-179609 (P2017-179609)	(71) 出願人	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
(22) 出願日	平成29年9月19日 (2017.9.19)	(74) 代理人	110002147 特許業務法人酒井国際特許事務所
		(72) 発明者	後藤 啓介 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
		(72) 発明者	丸橋 弘治 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
		(72) 発明者	稲越 宏弥 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

(54) 【発明の名称】 学習データ選択プログラム、学習データ選択方法、および、学習データ選択装置

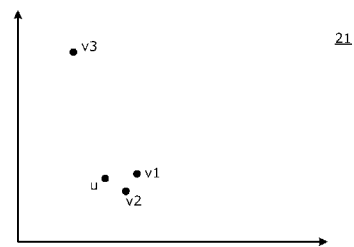
(57) 【要約】

【課題】 変換された入力データに対する分類・判別の要因を推定する機械学習モデルの、推定精度を向上させる。

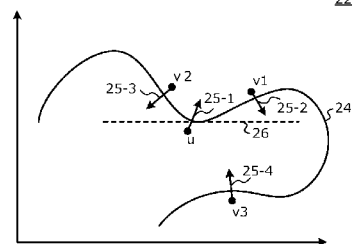
【解決手段】

入力データを変換した変換データに対し分類・判別を行う機械学習モデルの出力要因を推定する、推定モデルの学習データの選択を行うために、1) 機械学習モデルへの入力データ群に含まれる第1の入力データの指定に応じ、第1の入力データに関連する第1の入力データ群を抽出し、2) 第1の入力データ群に対応した、機械学習モデルに入力される第1の変換データ群、および、第1の変換データ群に対応した機械学習モデルの第1の出力データ群を、それぞれ取得し、3) 第1の入力データと第1の入力データ群のデータそれぞれとの距離、および、第1の変換データと第1の変換データ群のデータそれぞれとの距離に基づき、第1の入力データ群から、推定モデルの学習対象データを選択する。

入力データの特徴空間



変換データの特徴空間



【選択図】 図3

【特許請求の範囲】**【請求項 1】**

入力データを変換した変換データに対し分類または判別を行う機械学習モデルの出力要因を推定する、推定モデルの学習データの選択を、コンピュータに実行させる学習データ選択プログラムであって、

前記機械学習モデルへの入力データ群に含まれる第 1 の入力データの指定に応じ、前記第 1 の入力データに関連する第 1 の入力データ群を抽出し、

前記第 1 の入力データ群に対応した、前記機械学習モデルに入力される第 1 の変換データ群、および、前記第 1 の変換データ群に対応した前記機械学習モデルの第 1 の出力データ群を、それぞれ取得し、

前記第 1 の入力データと前記第 1 の入力データ群のデータそれぞれとの距離、および、前記第 1 の変換データと前記第 1 の変換データ群のデータそれぞれとの距離に基づき、前記第 1 の入力データ群から、前記推定モデルの学習対象データを選択する、

ことを特徴とする学習データ選択プログラム。

10

【請求項 2】

請求項 1 記載の学習データ選択プログラムであって、

前記抽出する処理は、前記第 1 の出力データ群に含まれる出力データそれぞれのデータ内容に基づき、学習対象データを抽出する、

ことを特徴とする学習データ選択プログラム。

20

【請求項 3】

請求項 2 記載の学習データ選択プログラムであって、

前記抽出する処理は、前記第 1 の出力データ群に含まれるデータ内容の比率に基づき、学習対象データを抽出する、

ことを特徴とする学習データ選択プログラム。

【請求項 4】

請求項 3 記載の学習データ選択プログラムであって、

前記抽出する処理は、前記第 1 の出力データ群に含まれるデータ内容の正例と負例の比率に基づき、学習対象データを抽出する、

ことを特徴とする学習データ選択プログラム。

30

【請求項 5】

請求項 1 記載の学習データ選択プログラムであって、

前記特定する処理は、データ取得タイミングが前記第 1 の入力データと所定の関係を有する入力データを前記第 1 の入力データ群と特定する、

ことを特徴とする学習データ選択プログラム。

【請求項 6】

請求項 1 記載の学習データ選択プログラムであって、

前記特定する処理は、データ生成元が前記第 1 の入力データと所定の関係を有する入力データを前記第 1 の入力データ群と特定する、

ことを特徴とする学習データ選択プログラム。

40

【請求項 7】

請求項 1 記載の学習データ選択プログラムであって、

前記抽出する処理は、前記第 1 の入力データと前記第 1 の入力データ群のデータそれぞれとの距離を、前記第 1 の入力データと前記第 1 の入力データ群のデータそれぞれとを個別に変換して算出する、

50

ことを特徴とする学習データ選択プログラム。

【請求項 8】

入力データを変換した変換データに対し分類または判別を行う機械学習モデルの出力要因を推定する、推定モデルの学習データの選択方法であって、

前記機械学習モデルへの入力データ群に含まれる第 1 の入力データの指定に応じ、前記第 1 の入力データに関連する第 1 の入力データ群を抽出し、

前記第 1 の入力データ群に対応した、前記機械学習モデルに入力される第 1 の変換データ群、および、前記第 1 の変換データ群に対応した前記機械学習モデルの第 1 の出力データ群を、それぞれ記憶装置より取得し、

前記第 1 の入力データと前記第 1 の入力データ群のデータそれぞれとの距離、および、前記第 1 の変換データと前記第 1 の変換データ群のデータそれぞれとの距離に基づき、前記第 1 の入力データ群から、前記推定モデルの学習対象データを選択する、

ことを特徴とする学習データ選択方法。

10

【請求項 9】

入力データを変換した変換データに対し分類または判別を行う機械学習モデルの出力要因を推定する、推定モデルの学習データの選択装置であって、

前記機械学習モデルへの入力データ群に含まれる第 1 の入力データの指定に応じ、前記第 1 の入力データに関連する第 1 の入力データ群を抽出する抽出部と、

前記第 1 の入力データ群に対応した、前記機械学習モデルに入力される第 1 の変換データ群、および、前記第 1 の変換データ群に対応した前記機械学習モデルの第 1 の出力データ群を、それぞれ記憶装置より取得する選択部と、を有し、

前記選択部は、前記第 1 の入力データと前記第 1 の入力データ群のデータそれぞれとの距離、および、前記第 1 の変換データと前記第 1 の変換データ群のデータそれぞれとの距離に基づき、前記第 1 の入力データ群から、前記推定モデルの学習対象データを選択する、

ことを特徴とする学習データ選択装置。

20

【発明の詳細な説明】

30

【技術分野】

【0001】

本発明は、機械学習の学習対象データを選択する、学習データ選択プログラム、学習データ選択方法、および、学習データ選択装置に関する。

【背景技術】

【0002】

近年、さまざまな分野のデータに対して、機械学習を用いた分類や判別が行われている。機械学習を用いることにより、精度の高い分類や判別が可能である一方、機械学習のどの特徴量が分類や判別の要因であるかは、一般的には知ることができない。

【0003】

機械学習が分類や判別をする際に、どのような要因により分類や判別を行ったのかが明確でない場合、たとえば、機械学習の適用分野を拡大する際のリスクとなる。

40

【0004】

機械学習に用いられる分類器による分類結果から、分類の要因となる特徴を推定する方法として、LIME (Local Interpretable Model-agnostic Explanations) という手法が知られている (例えば非特許文献 1 参照)。LIME においては、任意の分類器 f と入力データ u について、 u の分類結果 $f(u)$ に大きく貢献した u の要素・特徴を推定し、提示することが行われる。

【0005】

非特許文献 1 に記載された LIME の手法では、画像を対象とした分類結果に対する原

50

因推定が行われており、どの画像の部分が分類結果に寄与したかを推定することが記載されている。

【0006】

図2は、分類器による入力データの分類と、分類要因の推定の概要を示す図である。一般的な機械学習の分類器を用いた分類では、入力データを分類器が処理可能な次元に変換した変換データを生成し、生成された変換データに対して学習および分類が行われる。

【0007】

図2に示される一般的な機械学習の分類器を用いた分類に対し、LIMEの手法により分類結果に対する分類要因の推定を行うと、分類結果が変換データのどの要因に関連するかを推定するものとなり、入力データのどの要因に関連するものかを推定するものではないため、必ずしも有用であるとは限らなかった。

10

【先行技術文献】

【非特許文献】

【0008】

【非特許文献1】“Why Should I Trust You?” Explaining the Predictions of Any Classifier, Marco Tulio Ribeiro et. al., Proceedings of the 22nd ACM SIGKDD International Conference on Knowledge Discovery and Data Mining, 2016年8月

【発明の概要】

【発明が解決しようとする課題】

【0009】

上述したように、入力データを分類器が処理可能な次元に変換した変換データを生成し、生成された変換データに対して学習および分類が行われる機械学習の分類器では、入力データにおける分類要因の推定が求められる。

20

【0010】

図3は、入力データから生成された変換データの分類を行う機械学習の分類器に関する、各特徴空間における分類器への入力データと変換データの概要を示す図である。図3において、点uは正例として判別されるデータ、点v1、点v2、点v3は負例として判別されるデータに対応する、入力データの特徴空間21および変換データの特徴空間22での点である。

【0011】

特徴空間22において、点u、点v1、点v2、点v3にはそれぞれ、予測値が大きく変動する方向を表すベクトルである説明ベクトル25-1、25-2、25-3、25-4がそれぞれ示されている。また、変換データの特徴空間22において、近似識別線24は正例と負例の識別境界を近似する境界線である。入力データの特徴空間21では必ずしも分類器の正例と負例の識別境界が明確とは限らないため、図3の入力データの特徴空間21では近似識別線は描かれてはいない。

30

【0012】

ここで、図2に示した分類要因の推定を図3の点uの近傍について行うことは、正例と負例とを識別する近似識別線24の点uの近傍における識別要因を示す識別要因26を推定することに相当する。

40

【0013】

データuの近傍における分類要因を推定する場合、入力データの特徴空間21では、点u、点v1、点v2と点v3とは離れているため、点v3をLIMEの手法における点uの近傍の点として扱うことは適切ではない。

【0014】

非特許文献1に記載されたLIMEの手法により分類要因を推定する場合、変換データの特徴空間22内において点uの近傍の点を選択し、分類要因を推定する。具体的には、変換データの特徴空間22で点uの近傍に存在する点v1、点v2、点v3それぞれを学習データとして機械学習した推定器により、分類要因の推定を行う。すなわち、点v3という、変換データの特徴空間22においては点uの近傍にあるものの、入力データの特徴

50

空間 2 1 においては点 u の近傍にあるとはいえない点を学習データに含めて学習が行われるため、学習された推定モデルによる判別精度が劣化するという問題がある。

【 0 0 1 5 】

1 つの側面では、本件は、変換された入力データに対する分類または判別の要因を推定する機械学習モデルの、推定精度を向上させることを目的とする。

【課題を解決するための手段】

【 0 0 1 6 】

1 つの案では、入力データを変換した変換データに対し分類または判別を行う機械学習モデルの出力要因を推定する、推定モデルの学習データの選択について、コンピュータに以下の処理を実行させる学習データ選択プログラム、学習データ選択方法、および、学習データ選択装置が提供される。

10

【 0 0 1 7 】

すなわち、コンピュータに、機械学習モデルへの入力データ群に含まれる第 1 の入力データの指定に応じ、第 1 の入力データに関連する第 1 の入力データ群を抽出し、第 1 の入力データ群に対応した、機械学習モデルに入力される第 1 の変換データ群、および、第 1 の変換データ群に対応した機械学習モデルの第 1 の出力データ群を、それぞれ取得し、

第 1 の入力データと第 1 の入力データ群のデータそれぞれとの距離、および、第 1 の変換データと第 1 の変換データ群のデータそれぞれとの距離に基づき、第 1 の入力データ群から、推定モデルの学習対象データを選択することを特徴とする。

20

【発明の効果】

【 0 0 1 8 】

1 態様によれば、出力要因の推定精度を向上させる学習データを選択することができる。

【図面の簡単な説明】

【 0 0 1 9 】

【図 1】第 1 の実施形態に係る分類装置および要因推定装置の構成例を示す図である。

【図 2】分類器による入力データの分類と、分類要因の推定の概要を示す図である。

【図 3】各特徴空間における分類器への入力データと変換データの概要を示す図である。

【図 4】第 1 の実施形態に係る処理の手順の一例を示すフローチャートである。

30

【図 5】第 1 の実施形態のシステム構成例を示す図である。

【図 6】第 1 の実施形態に用いる監視サーバのハードウェアの一構成例を示す図である。

【図 7】第 2 の実施形態に係る要因推定装置の構成例を示す図である。

【図 8 A】選択データ u についての入力データと中間データとの関係を示す図である。

【図 8 B】対象データ $v_1 \sim v_3$ についての入力データと中間データとの関係を示す図である。

【図 9 A】選択データ u についての入力データと変換データとの関係を示す図である。

【図 9 B】対象データ $v_1 \sim v_3$ についての入力データと変換データとの関係を示す図である。

【図 1 0】選択データ u と対象データ $v_1 \sim v_3$ の各特徴空間における距離を示す図である。

40

【図 1 1】第 2 の実施形態における入力データ記憶部の一例を示す図である。

【図 1 2】第 2 の実施形態における変換データ記憶部の一例を示す図である。

【図 1 3】第 2 の実施形態における解析結果記憶部の一例を示す図である。

【図 1 4】第 2 の実施形態における入力データ距離記憶部の一例を示す図である。

【図 1 5】第 2 の実施形態における変換データ距離記憶部の一例を示す図である。

【図 1 6】第 2 の実施形態に係る処理の手順の一例を示すフローチャートである。

【図 1 7】第 2 の実施形態に係る入力データ距離算出部 2 2 1 の距離算出の一例を示す図である。

【図 1 8】変換行列を用いて生成した中間データ間の類似度の計算例を示す図である。

50

- 【図 19】「項 S」の変換行列の更新例を示す図である。
- 【図 20】類似度計算処理の手順の一例を示すフローチャートである。
- 【図 21】初期状態の変換行列を用いた中間データの生成例を示す図である。
- 【図 22】「項 S」の変換行列の更新例を示す図である。
- 【図 23】「項 R」の変換行列の更新例を示す図である。
- 【図 24】更新後の変換行列を用いた中間データの生成例を示す図である。
- 【図 25】類似度計算の比較例を示す第 1 の図である。
- 【図 26】類似度計算の比較例を示す第 2 の図である。
- 【図 27】類似度計算例を示す第 1 の図である。
- 【図 28】類似度計算例を示す第 2 の図である。

10

【発明を実施するための形態】

【0020】

以下、本実施の形態について図面を参照して説明する。なお各実施の形態は、矛盾のない範囲で複数の実施の形態を組み合わせることができる。

【0021】

〔第 1 の実施形態〕

図 1 は、第 1 の実施形態に係る分類装置 10 および要因推定装置 20 の構成例を示す図である。分類装置 10 は、収集部 110、入力データ記憶部 112、変換部 120、変換データ記憶部 122、分類部 130、学習結果記憶部 132、解析部 140、解析結果記憶部 142 を有する。

20

【0022】

要因推定装置 20 は、抽出部 210、選別部 220、選別データ記憶部 230、推定部 240、学習結果記憶部 242 を有する。例えばコンピュータが、学習データ選別プログラムを実行することによって、学習データの選別方法を実行可能な要因推定装置 20 が実現される。

【0023】

分類装置 10 において、収集部 110 は、入力データを収集し、入力データ記憶部 112 に記憶する。たとえば、図 3 に示される入力データの特徴空間 21 の入力データ u 、 v_1 、 v_2 、 v_3 を含む入力データが入力データ記憶部 112 に記憶される。変換部 120 は、収集部 110 で収集された入力データを所定の方法により変換し、変換データ記憶部 122 に記憶する。たとえば、図 3 に示される入力データの特徴空間 21 の入力データ u 、 v_1 、 v_2 、 v_3 は、変換データの特徴空間 22 の u 、 v_1 、 v_2 、 v_3 にそれぞれ変換される。所定の方法とは、たとえば、分類部 130 における分類が適切に行われるように入力データの次元、値等を変換するものであり、線形変換、非線形変換、一方向変換、双方向変換などを用いることができる。分類部 130 は、変換部 120 で変換された変換データを入力し、学習結果記憶部 132 に記憶された学習パラメータを用いて分類を行い、解析部 140 に出力する。

30

【0024】

要因推定装置 20 において、抽出部 210 は、たとえばユーザから、要因推定対象となる入力データの指定を受け、分類装置 10 の入力データ記憶部 112 から関連するデータを取得する。選別部 220 は、抽出部 210 で取得されたデータについて、分類装置 10 の変換データ記憶部 122 から変換データを、分類装置 10 の解析結果記憶部 142 から変換データに対応した解析結果を取得し、推定部 240 の学習対象となる学習データを選別する。推定部 240 は、選別部 220 により選別された学習データに基づき学習を行い、学習結果記憶部 242 に学習パラメータを記憶する。

40

【0025】

図 4 は、第 1 の実施形態に係る処理の手順の一例を示すフローチャートである。図 1 に示された要因推定装置 20 の抽出部 210 は、たとえばユーザから要因推定対象となる入力データである選択データとして、図 3 に示される点 u に対応した入力データ u の指定を受ける。入力データ u の指定を受けたことに対応して、抽出部 210 は、所定の基準によ

50

り分類装置 10 の入力データ記憶部 112 より、図 3 に示された点 $v_1 \sim v_3$ に対応した対象データである入力データ $v_1 \sim v_3$ を抽出する (S12)。所定の基準としては、例えば、入力データの収集タイミングが選択データ u と所定の関係を有するもの、具体的には、データ収集期間が選択データ u と前後 1 時間のもの、とすることができる。また、所定の基準として、例えば、データ生成元が選択データ u と所定の関係を有するもの、具体的には、同じサーバから対象データが取得されたもの、とすることもできる。

【0026】

次に、抽出部 210 からの選択データおよび対象データの抽出を受け、選別部 220 は、選択データ u および対象データ $v_1 \sim v_3$ に対応する変換データおよび分類結果を、分類装置 10 の変換データ記憶部 122 および解析結果記憶部 142 より読み込む (S14)。

10

【0027】

次に、選別部 220 は、選択データ u と対象データ $v_1 \sim v_3$ それぞれとの、入力データの特徴空間における距離と、変換データの特徴空間における距離とを算出し、距離に基づいて学習対象のデータを選別する (S16)。

【0028】

図 3 に示された選択データ u 、および、対象データ $v_1 \sim v_3$ は、入力データの特徴空間 21 においては、点 u と点 v_1 、点 v_2 との距離に比べ、点 u と点 v_3 との距離は大きい。一方、変換データの特徴空間 22 においては、点 u と点 $v_1 \sim v_3$ との距離は大きくは変わらない。

20

【0029】

すなわち、変換データの特徴空間 22 では点 u と近傍にある点 $v_1 \sim v_3$ のうち、点 v_3 は入力データの特徴空間 21 では点 u の近傍とはならない。このため、対象データ v_3 は、分類要因を推定する推定部 240 の学習データとしては不適切である。

【0030】

入力データの特徴空間 21 における点 u に対する距離と、変換データの特徴空間 22 における点 u に対する距離の和を取ると、対象データ v_3 の距離の和は、対象データ v_1 の距離の和、または、選択データ v_2 の距離の和よりも大きくなるから、距離の和の計算結果に基づき、対象データ v_3 を除外し、対象データ v_1 および v_2 を学習対象として選択することができる。

30

【0031】

次に、選別されたデータを用いて、推定部 240 において学習が行われる (S18)。すなわち、推定部 240 において、選択データ u と対象データ v_1 および v_2 の入力データ、変換データ、および、分類結果に基づく学習が行われ、学習結果に対応する学習パラメータが、学習結果記憶部 242 に記憶される。

【0032】

図 5 は、第 1 の実施形態のシステム構成例を示す図である。ネットワーク 2000 には、複数のサーバ 2011, 2012, …、複数の端末装置 2021, 2022, …、および監視サーバ 1001 が接続されている。複数のサーバ 2011, 2012, …は、いずれかの端末装置からの要求に応じた処理を実行するコンピュータである。複数のサーバ 2011, 2012, …のうちの 2 台以上が連携して処理を実行する場合もある。複数の端末装置 2021, 2022, …は、複数のサーバ 2011, 2012, …で提供されるサービスを利用するユーザが使用するコンピュータである。

40

【0033】

監視サーバ 1001 は、ネットワーク 2000 を介して行われた通信を監視し、通信ログを記録する。監視サーバ 1001 は、単位時間帯ごとの通信ログのデータを分類する。例えば監視サーバ 1001 は、対応する単位時間帯における不正通信の有無に応じて、データを分類する。

【0034】

監視サーバ 1001 上で動作する分類装置 10 の収集部 110 は、ネットワーク 200

50

0を介して送受信されているパケットなどの通信情報を取得する。例えば収集部110は、ネットワーク2000内に設置されたスイッチのミラーリングポートを介して、そのスイッチを経由して通信されたパケットを取得する。また収集部110は、各サーバ2011, 2012, …から、そのサーバ自身の通信ログを取得することもできる。収集部110は、取得した通信情報のログ(通信ログ)を、入力データ記憶部112に格納する。

【0035】

なお、図5に示した各要素間を接続する線は通信経路の一部を示すものであり、図示した通信経路以外の通信経路も設定可能である。また、図5に示した各要素の機能は、例えば、その要素に対応するプログラムモジュールをコンピュータに実行させることで実現することができる。

10

【0036】

図6は、第1の実施形態に用いる監視サーバのハードウェアの一構成例を示す図である。監視サーバ1010は、プロセッサ1011によって装置全体が制御されている。プロセッサ1011には、バス1019を介してメモリ1012と複数の周辺機器が接続されている。プロセッサ1011は、マルチプロセッサであってもよい。プロセッサ1011は、例えばCPU(Central Processing Unit)、MPU(Micro Processing Unit)、またはDSP(Digital Signal Processor)である。プロセッサ1011がプログラムを実行することで実現する機能の少なくとも一部を、ASIC(Application Specific Integrated Circuit)、PLD(Programmable Logic Device)などの電子回路で実現してもよい。

20

【0037】

メモリ1012は、監視サーバ1010の主記憶装置として使用される。メモリ1012には、プロセッサ1011に実行させるOS(Operating System)のプログラムやアプリケーションプログラムの少なくとも一部が一時的に格納される。また、メモリ1012には、プロセッサ1011による処理に必要な各種データが格納される。メモリ1012としては、例えばRAM(Random Access Memory)などの揮発性の半導体記憶装置が使用される。

【0038】

バス1019に接続されている周辺機器としては、ストレージ装置1013、グラフィック処理装置1014、入力インタフェース1015、光学ドライブ装置1016、機器接続インタフェース1017およびネットワークインタフェース1018がある。

30

【0039】

ストレージ装置1013は、内蔵した記録媒体に対して、電氣的または磁氣的にデータの書き込みおよび読み出しを行う。ストレージ装置1013は、コンピュータの補助記憶装置として使用される。ストレージ装置1013には、OSのプログラム、アプリケーションプログラム、および各種データが格納される。なお、ストレージ装置1013としては、例えばHDD(Hard Disk Drive)やSSD(Solid State Drive)を使用することができる。

【0040】

グラフィック処理装置1014には、モニタ1021が接続されている。グラフィック処理装置1014は、プロセッサ1011からの命令に従って、画像をモニタ1021の画面に表示させる。モニタ1021としては、CRT(Cathode Ray Tube)を用いた表示装置や液晶表示装置などがある。

40

【0041】

入力インタフェース1015には、キーボード1022とマウス1023とが接続されている。入力インタフェース1015は、キーボード1022やマウス1023から送られてくる信号をプロセッサ1011に送信する。なお、マウス1023は、ポインティングデバイスの一例であり、他のポインティングデバイスを使用することもできる。他のポインティングデバイスとしては、タッチパネル、タブレット、タッチパッド、トラックボ

50

ールなどがある。

【0042】

光学ドライブ装置1016は、レーザ光などを利用して、光ディスク1024に記録されたデータの読み取りを行う。光ディスク1024は、光の反射によって読み取り可能なようにデータが記録された可搬型の記録媒体である。光ディスク1024には、DVD (Digital Versatile Disc)、DVD-RAM、CD-ROM (Compact Disc Read Only Memory)、CD-R (Recordable) / RW (ReWritable) などがある。

【0043】

機器接続インタフェース1017は、監視サーバ1010に周辺機器を接続するための通信インタフェースである。例えば機器接続インタフェース1017には、メモリ装置1025やメモリアライタ1026を接続することができる。メモリ装置1025は、機器接続インタフェース1017との通信機能を搭載した記録媒体である。メモリアライタ1026は、メモリカード1027へのデータの書き込み、またはメモリカード1027からのデータの読み出しを行う装置である。メモリカード1027は、カード型の記録媒体である。

10

【0044】

ネットワークインタフェース1018は、ネットワーク1020に接続されている。ネットワークインタフェース1018は、ネットワーク1020を介して、他のコンピュータまたは通信機器との間でデータの送受信を行う。

【0045】

以上のようなハードウェア構成によって、第1の実施形態の処理機能を実現することができる。

20

【0046】

監視サーバ1010は、例えばコンピュータが読み取り可能な記録媒体に記録されたプログラムを実行することにより、第2の実施形態の処理機能を実現する。監視サーバ1010に実行させる処理内容を記述したプログラムは、様々な記録媒体に記録しておくことができる。例えば、監視サーバ1010に実行させるプログラムをストレージ装置1013に格納しておくことができる。プロセッサ1011は、ストレージ装置1013内のプログラムの少なくとも一部をメモリ1012にロードし、プログラムを実行する。また監視サーバ1010に実行させるプログラムを、光ディスク1024、メモリ装置1025、メモリカード1027などの可搬型記録媒体に記録しておくこともできる。可搬型記録媒体に格納されたプログラムは、例えばプロセッサ1011からの制御により、ストレージ装置1013にインストールされた後、実行可能となる。またプロセッサ1011が、可搬型記録媒体から直接プログラムを読み出して実行することもできる。

30

【0047】

以上のようなハードウェア構成によって、第1の実施形態の処理機能を実現することができる。

【0048】

〔第2の実施形態〕

次に第2の実施形態について説明する。

40

【0049】

図7は、第2の実施形態に係る要因推定装置の構成例を示す図である。図7に示された第2の実施形態に係る要因推定装置20のうち、第1の実施形態に係る要因推定装置10と同じ動作の部分については説明を省略する。

【0050】

第2の実施形態に係る要因推定装置20の選別部220は、入力データ距離算出部221、入力データ距離記憶部222、変換データ距離算出部223、変換データ距離記憶部224、対象判定部225を有する。

【0051】

図11は、第2の実施形態における、入力データ記憶部の一例を示す図である。図1に

50

示された分類装置 10 に含まれる入力データ記憶部 112 には、複数の単位期間ログ 112-1, 112-2, ... が格納されている。単位期間ログ 112-1, 112-2, ... それぞれには、通信ログの収集期間が示されており、例えば、単位期間ログ 112-1 の収集期間は、10:00 - 10:10、単位期間ログ 112-2 の収集期間は 10:10 - 10:20 である。単位期間ログ 112-1, 112-2, ... には、収集期間で示された時間帯内に収集した通信情報が格納される。

【0052】

単位期間ログ 112-1, 112-2, ... に格納される各レコードには、通信元ホスト、通信先ホスト、ポート、および量が含まれる。通信元ホストは、パケットの送信元の装置の識別子である。通信先ホストは、パケットの宛先の装置の識別子である。ポートは、通信元ホストと通信先ホストが通信を行った通信ポート番号の識別子である。単位期間ログ 112-1 における量は、通信元ホスト・通信先ホスト・ポートの組み合わせに対する値であり、たとえば、通信元ホスト・通信先ホスト・ポートの組が同じ通信の出現回数である。

10

【0053】

図 12 は、第 2 の実施形態における、変換データ記憶部の一例を示す図である。図 1 に示された分類装置 10 に含まれる変換データ記憶部 122 は、変換データテーブル 122-1, 122-2, ... を記憶している。各変換テーブル 122-1, 122-2, ... は、入力データ記憶部 112 に記憶された単位期間ログ 112-1, 112-2, ... を、分類部 130 における学習および分類に対応した変換により変換したデータである。

20

【0054】

入力データ記憶部 112 に記憶された単位期間ログ 112-1, 112-2, ... は、図 12 の各変換データの下部に示された順序づけ 127-1, 127-2, ... により、変換データテーブル 122-1, 122-2, ... に変換されている。

【0055】

分類装置 10 の学習段階において、分類部 130 の学習パラメータと各順序付けは、学習結果に応じた値および関係となっている。

【0056】

図 13 は、第 2 の実施形態における、解析結果記憶部の一例を示す図である。図 1 に示された分類装置 10 に含まれる解析結果記憶部 142 は、入力データ記憶部 112 に記憶された、収集期間の異なる単位期間ログ 112-1, 112-2, ... ごとの、分類部 130 による分類結果に対応した解析結果を記憶している。例えば、図 11 の入力データ記憶部 112 に記憶された、収集期間が 10:00 - 10:10 である単位期間ログ 112-1 は、分類部 130 による分類結果に基づく解析結果が、1 (問題あり) として、図 13 に示される解析結果テーブル 142-1 に保存されている。

30

【0057】

図 16 は、第 2 の実施形態に係る処理の手順の一例を示すフローチャートである。以下、第 2 の実施形態における、要因推定装置 20 の選別部 220 による、学習対象データの選別手順を、図 16 に示されたフローチャートと、図 8A、図 8B、図 9A、図 9B、および、図 10 を用いて説明する。

40

【0058】

図 7 に示された要因推定装置 20 の抽出部 210 は、たとえばユーザから要因推定対象となる入力データである選択データ u として、図 11 に示された収集期間が 10:00 - 10:10 である単位期間ログ 112-1 の指定を受ける。

【0059】

選択データ u の指定に対応して、抽出部 210 は、収集期間の始期が選択データ u に続く 10:10 - 10:30 である、図 11 に示された単位期間ログ 112-2 ~ 112-4 を対象データ v1 ~ v3 として選択する (S22)。

【0060】

50

ここでは、対象データの選定方法として、データ収集期間が選択データuと連続するものを選択しているが、実施例1と同様に、データ収集期間に関する他の基準や、データ生成元が選択データuと所定の関係を有するものを対象データとして選択してもよい。具体的には、同じサーバから対象データが取得されたものを対象データとして選択してもよい。

【0061】

次に、抽出部210により抽出された選択データuおよび対象データv1~v3について、入力データ距離算出部221は、選択データおよび各対象データを変換する順序付けを生成する(S24)。図8A、および、図8Bは、選択データuおよび対象データv1~v3についての、入力データと中間データとの関係を示す図である。図8Aおよび図8Bにおいて、選択データuの入力データ801、および、対象データv1~v3の入力データ811~831は、順序付け803~833により中間データ805~835に変換される。順序付け803~833の算出方法については、図17~図28を参照し後述する。

10

【0062】

順序付け803~833を算出することにより、選択データuの入力データ801、および、対象データv1~v3の入力データ811~831は、中間データ805~835に変換され、入力データの特徴空間21における選択データuと対象データv1~v3の距離を、それぞれ求めることができる。

【0063】

たとえば、選択データuと対象データv1との、入力データの特徴空間21における距離は、中間データ805と中間データ815との距離を求めることにより算出される。

20

【0064】

中間データ805と中間データ815についての計算例を以下に示す。

a) uとv1の双方に存在する量が異なる項目は、(S'2、R'3、P'2)の1項目であるので、 $(2-1)^2=2$ 。

b) uのみに存在する項目は、(S'1、R'3、P'1)および(S'3、R'1、P'1)の2項目であり、 $1^2+1^2=2$ 。

c) v1のみに存在する項目は、(S'2、R'3、P'1)および(S'2、R'2、P'1)の2項目であり、 $1^2+1^2=2$ 。となるから、a)b)c)の合計は5となる。

30

【0065】

同様に計算することにより、入力データの特徴空間21における選択データuと対象データv1~v3の距離は、それぞれ、5、4、9となる。

【0066】

次に、抽出部210からの選択データおよび対象データの抽出を受け、選別部220は、選択データuおよび対象データv1~v3に対応する変換データおよび分類結果を、図1に示された分類装置10の変換データ記憶部122および解析結果記憶部142より読み込む(S26)。

【0067】

次に、入力データの特徴空間21におけるテンソル間距離および変換データの特徴空間22におけるテンソル間距離に基づき学習対象のデータを選別する(S28)。図9A、および、図9Bは、選択データuおよび対象データv1~v3についての、入力データと変換データとの関係を示す図である。図9Aおよび図9Bにおいて、選択データuの入力データ901、および、対象データv1~v3の入力データ911~931は、分類装置10の変換部120により生成された順序付け903~933により、変換データ905~935に変換されている。

40

【0068】

入力データの特徴空間21における距離算出と同様に、変換データの特徴空間22における選択データuと対象データv1~v3の距離を求める。選択データuの変換データ9

50

05と、対象データ $v_1 \sim v_3$ の変換データ915～935との距離は、それぞれ、9、8、9となる。

【0069】

図14は、第2の実施形態における、入力データ距離記憶部222の一例を示す図である。入力データ距離記憶部222には、入力データ距離算出部により算出された距離が、入力データ距離テーブル222-1として、収集期間毎に記憶される。

【0070】

図15は、第2の実施形態における、変換データ距離記憶部224の一例を示す図である。変換データ距離記憶部224には、変換データ距離算出部により算出された距離が、変換データ距離テーブル224-1として、収集期間毎に記憶される。

10

【0071】

図10は、選択データ u と対象データ $v_1 \sim v_3$ の各特徴空間における距離を示す図である。図10において、選択データ u に対する、対象データ $v_1 \sim v_3$ の各特徴空間における距離の和は、それぞれ、14、12、18となる。

【0072】

ここで、図10および図3を参照すると、選択データ u 、および、対象データ $v_1 \sim v_3$ は、入力データの特徴空間21においては、点 u と点 v_1 、点 v_2 との距離(4)に比べ、点 u と点 v_3 との距離(9)は大きい。一方、変換データの特徴空間22においては、点 u と点 $v_1 \sim v_3$ との距離は大きくは変わらない(9または8)。

【0073】

対象データ $v_1 \sim v_3$ について、選択データ u に対する各特徴空間の距離の和を取ると、上述したように、14、12、18となるから、対象データ v_3 は、対象データ v_1 および v_2 と比較して、選択データ u と離れており、分類要因を推定する推定部240の学習データとしては不適切である。距離の和の計算結果に基づき、対象データ v_3 を除外し、選択データ v_1 および v_2 を学習対象として選択することができる。

20

【0074】

図17は、第2の実施形態に係る入力データ距離算出部221の距離算出の一例を示す図である。入力データ距離算出部221は、抽出部210からの入力データについて算出した類似度に基づき、距離を算出する。

【0075】

以下では、簡単のために類似度を計算する第1データ1と第2データ2が、それぞれ2つの項目を持つ場合について説明する。項目数が3以上の場合であっても、以下の説明における行列をテンソルに拡張することで類似度を計算することが可能である。

30

【0076】

図17において、第1データ1と第2データ2は、入力データ距離算出部221による類似度の算出対象である。第1データ1は複数の第1レコードを有し、第1レコードの各々は、第1項目である「項S、項R」のそれぞれについての第1項目値(「項S」についての「 S_1, S_2 」、 「項R」についての「 R_1, R_2 」)と、第1項目値間の関係を示す数値「 $k_{11} \sim k_{14}$ 」を有する。同様に、第2データ2は複数の第2レコードを有し、第2レコードの各々は、第2項目である「項S、項R」のそれぞれについての第2項目値(「項S」についての「 S_1, S_2 」、 「項R」についての「 R_1, R_2 」)と、第2項目値間の関係を示す数値「 $k_{21} \sim k_{24}$ 」を有する。

40

【0077】

入力データ距離算出部221は、対象項目ごとに、第1データ1内の対象項目に属する対象第1項目値に関する他の第1項目値との関係と、第2データ2内の対象項目に属する対象第2項目値に関する他の第2項目値との関係との類似度を計算する。次に、計算した類似度に基づいて、第1重み情報5,6と第2重み情報7,8とを生成する。第1重み情報5,6は、複数の項目「項S、項R」のいずれかに属する複数の変換先項目値「 S'_1, S'_2, R'_1, R'_2 」のうちの、対象項目に属する対象変換先項目値への対象第1項目値の影響度を示す情報である。第2重み情報7,8は、対象変換先項目値への対象第

50

2項目値の影響度を示す情報である。

【0078】

例えば「項S」が対象項目として選択されたとき、入力データ距離算出部221は、第1重み情報5と第2重み情報7とを生成する。このとき第1重み情報5には、第1データ1内の「項S」に属する項目値「S1, S2」それぞれと、第2データ2内の「項S」に属する項目値「S1, S2」それぞれとの対ごとの、他の項目との関係の類似度が維持できるように、重みが設定される。

【0079】

同様に、第2重み情報7にも、第1データ1内の「項S」に属する項目値「S1, S2」それぞれと、第2データ2内の「項S」に属する項目値「S1, S2」それぞれとの対ごとの、他の項目との関係の類似度が維持できるように、重みが設定される。「項R」が対象項目として選択されたときに、入力データ距離算出部221は、第1重み情報6と第2重み情報8とを生成する。

10

【0080】

対象第1項目値と対象第2項目値との類似度を計算する場合、入力データ距離算出部221は、例えば複数の項目「項S, 項R」それぞれについて、初期値が設定された第1重み情報5, 6と第2重み情報7, 8とを生成する。

【0081】

次に、入力データ距離算出部221は、対象項目以外の項目について生成された他項目第1重み情報と他項目第2重み情報とに基づいて、第1データ1内の対象項目に属する対象第1項目値と、第2データ2内の対象項目に属する対象第2項目値との類似度を計算する。

20

【0082】

なお、入力データ距離算出部221は、所定の終了条件を満たすまで、複数の項目「項S, 項R」それぞれを、繰り返し対象項目として特定し、対象項目に対する第1重み情報5, 6と第2重み情報7, 8とを繰り返し生成してもよい。例えば、入力データ距離算出部221は、対象項目以外の項目について生成された他項目第1重み情報と他項目第2重み情報とを用いて、対象項目に属する対象第1項目値それぞれと対象第2項目値それぞれとの類似度を計算し、類似度に応じて対象項目の重み情報を更新する。

【0083】

具体的には、対象項目が「項S」であれば、入力データ距離算出部221は、「項R」について生成された第1重み情報6を用いて、第1データ1の項目値「S1, S2」それぞれと、第2データ2の項目値「S1, S2」それぞれとの類似度を計算する。次に、入力データ距離算出部221は、計算した類似度に基づいて、「項S」についての第1重み情報5と第2重み情報7とを更新する。さらに、入力データ距離算出部221は、更新後の第1重み情報5と第2重み情報7を用いて、「項R」についての第1重み情報6と第2重み情報8とを更新する。

30

【0084】

このように第1重み情報5, 6と第2重み情報7, 8の更新を繰り返すことで、第1類似判断用データ3と第2類似判断用データ4との類似度が向上するように、第1重み情報5, 6と第2重み情報7, 8が最適化される。

40

【0085】

次に、入力データ距離算出部221は、複数の項目「項S, 項R」それぞれについて生成された第1重み情報5, 6に基づいて、第1データ1を第1類似判断用データ3に変換する。第1類似判断用データ3は、複数の変換先項目値「S'1, S'2, R'1, R'2」のうちの異なる項目に属する2以上の変換先項目値間の関係を示す数値「k31~k34」が設定された複数の第3レコードを有する。さらに、入力データ距離算出部221は、複数の項目「項S, 項R」それぞれについて生成された第2重み情報7, 8に基づいて、第2データ2を第2類似判断用データ4に変換する。第2類似判断用データ4は、複数の変換先項目値「S'1, S'2, R'1, R'2」のうちの異なる項目に属する2以

50

上の変換先項目値間の関係を示す数値「 $k_{41} \sim k_{44}$ 」が設定された複数の第4レコードを有する。

【0086】

さらに、入力データ距離算出部221は、第1類似判断用データ3内の複数の第3レコードに含まれる数値群と、第2類似判断用データ4内の複数の第4レコードに含まれる数値群との類似度を計算する。入力データ距離算出部221は、第1重み情報5, 6と第2重み情報7, 8を繰り返し生成するとき、第1重み情報5, 6と第2重み情報7, 8を生成するごとに、第1類似判断用データ3の数値群と第2類似判断用データ4の数値群との類似度を計算する。そして、入力データ距離算出部221は、計算した類似度の最大値を、第1データ1と第2データ2との類似度と判定する。

10

【0087】

このように、第1重み情報5, 6と第2重み情報7, 8を用いて第1データ1と第2データ2とを変換した上で、類似度を計算することで、精度の高い類似度を算出することができる。すなわち、第1データ1と第2データ2との同一の項目に属する項目値に関する、他の項目値との間の関係の類似度が、その項目に対応する第1重み情報5, 6と第2重み情報7, 8とに反映されている。これにより、第1データ1と第2データ2との同一の項目に属する項目値に関する、他の項目値との間の関係の類似度が高いほど、変換後の第1類似判断用データ3と第2類似判断用データ4との類似度が高くなる。その結果、類似度の判定精度が向上し、第1データと第2データとの距離についても精度良く算出することができる。

20

【0088】

しかも、組み合わせ爆発のような計算量の急激な増加は発生せず、現実的な処理量での類似度計算が可能である。例えば、類似度の計算処理は、行列を用いて以下のように計算できる。

【0089】

入力データ距離算出部221は、特定の項目の項目値と他の項目との関係をベクトルで表現する。そして、入力データ距離算出部221は、第1データ1と第2データ2とを、2つの項目値に対応するベクトル間の距離を保持したまま、第1類似判断用データ3と第2類似判断用データ4とに変換する。このとき、入力データ距離算出部221は、変換に用いる第1重み情報5, 6と第2重み情報7, 8とを行列で表す。以下、第1重み情報5, 6と第2重み情報7, 8とを表す行列を、変換行列と呼ぶ。

30

【0090】

入力データ距離算出部221は、第1類似判断用データ3と第2類似判断用データ4間の最大類似度を、第1データ1と第2データ2との間の類似度とする。これにより、本質的な関係の構造に基づく類似度を計算でき、第1データ1と第2データ2との間の本質的な関係の構造に基づく距離を計算することができる。

【0091】

以下に、変換行列を用いた類似度の計算に関する詳細を説明する。上述のように、簡単のため、第1データ1と第2データ2との項目は2つだけとする。入力データ距離算出部221は、第1データ1と第2データ2とを、行列 X_1, X_2 で表す。行列 X_1, X_2 の各行は、1つ目の項目「項S」の各項目値「 S_1, S_2 」に対応し、各列は2つめの項目「項R」の各変数値「 R_1, R_2 」に対応する。行列の要素(成分)には、行に対する項目値と列に対応する項目値との関係を示す数値が入る。

40

【0092】

なお、入力データ距離算出部221は、1つ目の項目「項S」の項目値の種類数が第1データ1と第2データ2とで異なる場合には、少ないほうのデータにダミーの項目値を追加して、種類数を同数にする。入力データ距離算出部221は、2つ目の項目「項R」についても同様に、項目値の種類数を同数に揃える。

【0093】

入力データ距離算出部221は、第1データ1の「項S」と「項R」とに関する変換行

50

列（第1重み情報5, 6）を、それぞれ正方行列 C_{11} と C_{12} で表す。同様に、入力データ距離算出部221は、第2データ2の「項S」と「項R」とに関する変換行列（第2重み情報7, 8）を、それぞれ正方行列 C_{21} と C_{22} で表す。ただし、 C_{11} , C_{12} , C_{21} , C_{22} は、いずれも以下の正規直交条件を満たすものとする。

【0094】

【数1】

$$C_{11}^T C_{11} = C_{21}^T C_{21} = I$$

$$C_{12}^T C_{12} = C_{22}^T C_{22} = I$$

【0095】

I は対角成分が「1」で残りが「0」の単位行列である。このとき、 X_1 の列ベクトルを x_{1a} , x_{1b} とする。 x_{1a} , x_{1b} は、「項R」の変数値「a」、「b」と「項S」との関係を表しており、以下の関係を有する。

【0096】

【数2】

$$\|C_{11}^T x_{1a} - C_{11}^T x_{1b}\|^2 = \|x_{1a} - x_{1b}\|^2$$

【0097】

すなわち、 C_{11} による X_1 の変換は、項目値の他項目との関係を表すベクトル間の距離を変化させない。 C_{12} , C_{21} , C_{22} についても同様である。

【0098】

入力データ距離算出部221は、 C_{11} と C_{21} の更新では、 C_{12} と C_{22} を固定したときの、データ間類似度を最大化する C_{11} と C_{21} として算出する。データ間類似度 $E(X_1, X_2)$ は、以下の式で表される。

【0099】

【数3】

$$E(X_1, X_2) = \langle C_{11}^T X_1 C_{12}, C_{21}^T X_2 C_{22} \rangle$$

【0100】

データ間類似度を最大化する C_{11} と C_{21} は、以下に示す特異値分解により算出できる。

【0101】

【数4】

$$C_{11} S C_{21}^T = X_1 C_{12} C_{22}^T X_2^T$$

【0102】

ただし、 S は非負値を持つ正方対角行列である。

【0103】

このようにして、行列を用いて効率的にデータ変換を行い、類似度を計算することができる。

【0104】

このとき、行列 X_1 , X_2 でそれぞれ表される第1データ1と第2データ2との距離 $D(X_1, X_2)$ は、以下となる。

【0105】

【数5】

$$D(X_1, X_2) = \|C_{11}^T X_1 C_{12}, C_{21}^T X_2 C_{22}\|$$

【0106】

図17の例では、第1重み情報5, 6および第2重み情報7, 8が変換行列で表されている。例えば、第1重み情報5を示す変換行列の第1行・第1列の成分には、第1データ1の「項S」に属する項目値「S1」の、「項S」に属する変換先項目値「S'1」への

10

20

30

40

50

影響を示す重み (w_{11}) が設定されている。変換行列を用いると、第1データ1の項目ごとの項目値を成分とする行ベクトルに右から変換行列を乗算すれば、変換先項目値を得ることができる。例えば第1データ1の「項S」に属する項目値を成分とする行ベクトル (S_1, S_2) に、「項S」に関する第1重み情報5を示す変換行列を右から掛けることで、「項S」に属する変換先項目値を示す行ベクトル (S'_1, S'_2) が得られる。

【0107】

同様に、第1データ1の「項R」に属する項目値を成分とする行ベクトル (R_1, R_2) に、「項R」に関する第1重み情報6を示す変換行列を右から掛けることで、「項R」に属する変換先項目値を示す行ベクトル (R'_1, R'_2) が得られる。第2データ2の「項S」に属する項目値を成分とする行ベクトル (S_1, S_2) に、「項S」に関する第2重み情報7を示す変換行列を右から掛けることで、「項S」に属する変換先項目値を示す行ベクトル (S'_1, S'_2) が得られる。第2データ2の「項R」に属する項目値を成分とする行ベクトル (R_1, R_2) に、「項R」に関する第2重み情報8を示す変換行列を右から掛けることで、「項R」に属する変換先項目値を示す行ベクトル (R'_1, R'_2) が得られる。

10

【0108】

ここで第1データ1と第2データ2における同一レコード内の「項S」の項目値と「項R」の項目値の乗算結果が、そのレコードの「数値」の値であるものとする。同様に、第1類似判断用データ3と第2類似判断用データ4における同一レコード内の「項S」の項目値と「項R」の項目値の乗算結果が、そのレコードの「数値」の値であるものとする。すると、第1類似判断用データ3と第2類似判断用データ4との「数値」の値を算出できる。例えば第1類似判断用データ3の「 S'_1 」と「 R'_1 」との組に対応する数値「 k_{31} 」は、以下の通りとなる。

20

$$\begin{aligned} k_{31} &= S'_1 \times R'_1 \\ &= (w_{11} \times S_1 + w_{12} \times S_2) \times (w_{21} \times R_1 + w_{22} \times R_2) \\ &= w_{11} \times w_{21} \times S_1 \times R_1 + w_{12} \times w_{21} \times S_2 \times R_1 + w_{11} \times w_{22} \times S_1 \times R_2 + w_{12} \times w_{22} \times S_2 \times R_2 \\ &= w_{11} \times w_{21} \times k_{11} + w_{12} \times w_{21} \times k_{12} + w_{11} \times w_{22} \times k_{13} + w_{12} \times w_{22} \times k_{14} \end{aligned}$$

30

同様にして、第1類似判断用データ3と第2類似判断用データ4との「数値」の他の値 ($k_{31} \sim k_{34}, k_{41} \sim k_{44}$) も算出できる。

【0109】

入力データ距離算出部221は、第1類似判断用データ3と第2類似判断用データ4の「数値」の各値を比較することで、第1類似判断用データ3と第2類似判断用データ4との類似度を計算する。例えば、入力データ距離算出部221は、第1類似判断用データ3の各レコードの数値を成分とするベクトルと、第1類似判断用データ4の各レコードの数値を成分とするベクトルとの内積を計算し、内積の結果を類似度とする。

【0110】

このように第1重み情報5, 6および第2重み情報7, 8を行列で表すことで、類似度を計算することができる。

40

【0111】

なお上記の計算例は、簡単のために第1データ1と第2データ2との項目は2つだけとしているが、項目数が多い場合、行列をテンソルに拡張することで類似度を計算できる。なお、行列は、テンソルの一例である。

【0112】

比較対象のデータに対応するテンソルを X_m, X_n とする (m, n はデータを識別する整数)。 X_m, X_n に含まれる項目数が k (k は2以上の整数) のとき、変換行列を C_k とすると、データの類似判断用データへの変換は、以下の式で表すことができる。

【0113】

50

【数 6】

$$x_n \prod_k x_k C_k \dots (5)$$

【0 1 1 4】

式(5)の x_k は、テンソルのモード積を示している。式(5)の結果を用いて、テンソルを X_m 、 X_n 間の距離を、以下の式で表すことができる。

【0 1 1 5】

【数 7】

$$E(X_m, X_n) = \|X_m\|_2 + \|X_n\|_2 - 2 \langle X_m, X_n \prod_k x_k C_k \rangle \dots (6)$$

10

【0 1 1 6】

式(6)に示す距離を最小にする行列 C_k が変換行列となる。ただし、 C_k は、以下の正規直交条件を満たすものとする。

【0 1 1 7】

【数 8】

$$\begin{cases} C_k^T C_k = I \ (I_{kn} \geq I_{km}) \\ C_k C_k^T = I \ (I_{kn} < I_{km}) \end{cases} \dots (7)$$

20

【0 1 1 8】

項目ごとの C_k は、以下の特異値分解を、項目ごとに交互に繰り返し行うことで算出できる。

【0 1 1 9】

【数 9】

$$P_k S_k Q_k^T = \left(X_n \prod_{\substack{k'=1, \dots, K \\ k' \neq k}} x_{k'} C_{k'} \right)^{(k)T} X_m^{(k)} \dots (8)$$

30

【0 1 2 0】

式(8)の(k)は、テンソルを、第k番目の項目を列、その他の項目を行とする行列に変換する操作を表す。式(8)により、行列 $P_k S_k Q_k^T$ が生成される。そして P_k と Q_k^T を用いて、以下の式により行列 C_k が得られる。

【0 1 2 1】

【数 10】

$$C_k = P_k Q_k^T \dots (9)$$

【0 1 2 2】

このような計算により変換行列を求めデータを変換することで、項目数が3以上であっても現実的な計算量で類似度を計算可能であり、距離を計算することもできる。

40

【0 1 2 3】

図18は、変換行列を用いて生成した中間データ間の類似度の計算例を示す図である。入力データ距離算出部221は、第1データ1031と第1データ1032それぞれに対して、「量」以外の変数の数に応じた変換行列1041～1044を生成する。例えば入力データ距離算出部221は、第1データ1031に対して、「項S」に対応する変換行列1041と「項R」に対応する変換行列1042とを生成する。同様に入力データ距離算出部221は、第2データ1032に対して、「項S」に対応する変換行列1043と「項R」に対応する変換行列1044とを生成する。

【0 1 2 4】

50

変換行列 1041 ~ 1044 は、正規直交条件を満たす 2 行 2 列の行列である。変換行列 1041 の各行には、第 1 データ 1031 における「項 S」の変数値「S1」、「S2」が関連付けられている。変換行列 1043 の各行には、第 2 データ 1032 における「項 S」の変数値「S1」、「S2」が関連付けられている。変換行列 1041 の各列には、中間データ 1051 における「項 S」の変数値「S'1」、「S'2」が関連付けられている。変換行列 1043 の各列には、中間データ 1052 における「項 S」の変数値「S'1」、「S'2」が関連付けられている。変換行列 1041, 1043 の各成分には、行方向に関連付けられた変数値「S1」、「S2」を、列方向に関連付けられた変数値「S'1」、「S'2」へ変換する場合の重みが設定されている。

【0125】

変換行列 1042 の各行には、第 1 データ 1031 における「項 R」の変数値「R1」、「R2」が関連付けられている。変換行列 1044 の各行には、第 2 データ 1032 における「項 R」の変数値「R1」、「R2」が関連付けられている。変換行列 1042 の各列には、中間データ 1051 における「項 R」の変数値「R'1」、「R'2」が関連付けられている。変換行列 1044 の各列には、中間データ 1052 における「項 R」の変数値「R'1」、「R'2」が関連付けられている。変換行列 1042, 1044 の各成分には、行方向に関連付けられた変数値「R1」、「R2」を、列方向に関連付けられた変数値「R'1」、「R'2」へ変換する場合の重みが設定されている。

【0126】

なお図 18 の例では、変換行列 1041 ~ 1044 に設定されている重みの値を小数点 2 桁までしか示していないが、実際には小数点 2 桁よりも下位の桁の値も存在するものとする。

【0127】

入力データ距離算出部 221 は、変換行列 1041, 1042 を用いて、第 1 データ 1031 を中間データ 1051 に変換する。中間データ 1051 には、「項 S」の変数値と「項 R」の変数値の組み合わせに対応する量が設定されている。中間データ 1051 の「項 S」の変数値には「S'1」または「S'2」が設定され、「項 R」の変数値には「R'1」または「R'2」が設定されている。

【0128】

中間データ 1051 の「量」の値は、「項 S」の変数値と「項 R」の変数値との乗算結果である。入力データ距離算出部 221 は、第 1 データ 1031 を変換行列 1041, 1042 で変換することで、中間データ 1051 の「量」の値を算出する。例えば変数値「S'1」は、重み「-0.68」×「S1」+重み「-0.73」×「S2」である。変数値「S'2」は、重み「-0.73」×「S1」+重み「0.68」×「S2」である。変数値「R'1」は、重み「-0.32」×「R1」+重み「-0.94」×「R2」である。変数値「R'2」は、重み「-0.94」×「R1」+重み「0.32」×「R2」である。

【0129】

このように、変換行列 1041, 1042 を用いて、変数値「S'1」、「S'2」、「R'1」、「R'2」の値を、変数値「S1」、「S2」、「R1」、「R2」と、それらの重みで表すことができる。すると、変数値「S'1」または「S'2」と変数値「R'1」または「R'2」とを乗算すると、「S1×R1」、「S2×R1」、「S1×R2」、「S2×R2」のいずれを含む項が現れる。例えば「S1×R1」は、以下の式で表される。

【0130】

$$\begin{aligned} S'1 \times R'1 &= \{ (-0.68 \times S1) + (-0.73 \times S2) \} \times \{ (-0.32 \times R1) + (-0.94 \times R2) \} \\ &= (-0.68) \times (-0.32) \times S1 \times R1 + (-0.73) \times (-0.32) \times S2 \times R1 \\ &\quad + (-0.68) \times (-0.94) \times S1 \times R2 + (-0.73) \times (-0.94) \times S2 \times R2 \end{aligned}$$

10

20

30

40

50

入力データ距離算出部 221 は、「 $S_1 \times R_1$ 」、「 $S_2 \times R_1$ 」、「 $S_1 \times R_2$ 」、「 $S_2 \times R_2$ 」の値として、第 1 データ 1031 における対応する「量」の値を代入する。図 18 の例では、「 $S_1 \times R_1 = 1$ 」、「 $S_2 \times R_1 = 0$ 」、「 $S_1 \times R_2 = 1$ 」、「 $S_2 \times R_2 = 1$ 」である。その結果、中間データ 1051 における「量」の値が求まる。同様に、入力データ距離算出部 221 は、第 2 データ 1032 を変換行列 1043, 1044 で変換して、中間データ 1052 を生成する。

【0131】

入力データ距離算出部 221 は、中間データ 1051, 1052 の間の類似度を計算する。例えば入力データ距離算出部 221 は、中間データ 1051 の「量」の各変数値を成分とするベクトルと、中間データ 1052 の「量」の各変数値を成分とするベクトルとを、長さ「1」に正規化後、内積を計算する。そして入力データ距離算出部 221 は、内積の結果を、中間データ 1051, 1052 間の類似度とする。

10

【0132】

このようにして計算される中間データ 1051, 1052 間の類似度は、変換行列 1041 ~ 1044 に設定されている重みに依存する。そこで入力データ距離算出部 221 は、類似度が高くなるように変換行列 1041 ~ 1044 を更新する。変換行列 1041 ~ 1044 の更新は、「項 S」の変換行列 1041, 1043 の更新と、「項 R」の変換行列 1042, 1044 の更新とが交互に行われる。

【0133】

図 19 は、「項 S」の変換行列の更新例を示す図である。「項 S」の変換行列 1041, 1043 を更新する場合、入力データ距離算出部 221 は、「項 S」の変数を固定とし、「項 S」以外の変数を変換して中間データ 1053, 1054 を生成する。図 19 の例では、入力データ距離算出部 221 は、「項 R」の変数値「 R_1 」「 R_2 」を変換行列 1042, 1044 を用いて変換し、中間データ 1053, 1054 を生成している。中間データ 1053, 1054 の「量」の値は、「 S_1 」または「 S_2 」と「 R'_1 」または「 R'_2 」との乗算結果である。例えば第 1 データ 1031 の中間データ 1053 における「 $S_1 \times R'_1$ 」は、変換行列 1042 に示される重みを用いて、「 $(-0.32) \times S_1 \times R_1 + (-0.94) \times S_1 \times R_2$ 」と表される。第 1 データ 1031 に基づいて、「 $S_1 \times R_1$ 」と「 $S_1 \times R_2$ 」とに値を設定すれば、「 $S_1 \times R'_1$ 」の値が得られる。

20

30

【0134】

第 1 データ 1031 と第 2 データ 1032 の中間データ 1053, 1054 が生成されると、入力データ距離算出部 221 は、中間データ 1053 における「 S_1 」、「 S_2 」それぞれと、中間データ 1054 における「 S_1 」、「 S_2 」それぞれとの類似度を計算し、類似度行列 1061 を生成する。類似度行列 1061 の各行には、第 1 データ 1031 の「項 S」の変数値が関連付けられており、類似度行列 1061 の各列には、第 1 データ 1032 の「項 S」の変数値が関連付けられている。類似度行列 1061 の成分には、その成分が設定された行の変数値と列の変数値との類似度が設定されている。

【0135】

例えば入力データ距離算出部 221 は、「項 S」の各変数値それぞれについて、他の「項 R」の各変数値との関係を示すベクトルを生成する。具体的には、入力データ距離算出部 221 は、中間データ 1053 の「 S_1 」について、「 R'_1 」と「 R'_2 」とのそれぞれとの関係を示す「量」の値を成分とするベクトル v_{1_1} を生成する。同様に入力データ距離算出部 221 は、中間データ 1053 の「 S_2 」について、ベクトル v_{2_1} を生成する。入力データ距離算出部 221 は、中間データ 1054 の「 S_1 」について、ベクトル v_{1_2} を生成する。入力データ距離算出部 221 は、中間データ 1054 の「 S_2 」について、ベクトル v_{2_2} を生成する。

40

【0136】

入力データ距離算出部 221 は、ベクトル v_{1_1} とベクトル v_{1_2} との内積を、第 1 デー

50

タ1031の「S1」と第2データ1032の「S1」との類似度として、類似度行列1061に設定する。入力データ距離算出部221は、ベクトル v_{1_1} とベクトル v_{2_2} との内積を、第1データ1031の「S1」と第2データ1032の「S2」との類似度として、類似度行列1061に設定する。入力データ距離算出部221は、ベクトル v_{2_1} とベクトル v_{1_2} との内積を、第1データ1031の「S2」と第2データ1032の「S1」との類似度として、類似度行列1061に設定する。入力データ距離算出部221は、ベクトル v_{2_1} とベクトル v_{2_2} との内積を、第1データ1031の「S2」と第2データ1032の「S2」との類似度として、類似度行列1061に設定する。

【0137】

入力データ距離算出部221は、このようにして生成した類似度行列1061に基づいて、第1データ1031の「項S」変換用の変換行列1041aと第1データ1032の「項S」変換用の変換行列1043aとを生成する。例えば入力データ距離算出部221は、変換行列1041a, 1043aから $S'1$ 、 $S'2$ を消去して1つの行列にしたときに類似度行列1061に最も類似するような、変換行列1041a, 1043aを生成する。具体的には、入力データ距離算出部221は、類似度行列1061を特異値分解し、変換行列1041a, 1043aを生成する。

10

【0138】

類似度行列1061は、第1データ1031の「項S」変換用の変換行列1041を、生成した変換行列1041aに更新する。また類似度行列1061は、第2データ1032の「項S」変換用の変換行列1043を、生成した変換行列1043aに更新する。

20

【0139】

このようにして、「項S」を固定して他の変数を変換することで、「項S」の変換行列が更新される。次に入力データ距離算出部221は、「項R」を固定して他の変数を変換することで、「項R」の変換行列を更新する。入力データ距離算出部221は、各変数の変換行列を更新したら、更新後の変換行列を用いて、第1データ1031と第1データ1032との中間データを生成し、中間データ間の類似度を計算する。入力データ距離算出部221は、例えば中間データ間の類似度が収束するまで、変換行列の更新を繰り返し行う。これにより中間データ間の類似度の最大値を得る変換行列が生成される。そして入力データ距離算出部221は、中間データ間の類似度の最大値を、第1データ1031と第2データ1032との類似度とする。

30

【0140】

図18、図19に示した処理の手順をフローチャートで表すと図20のようになる。

【0141】

図20は、類似度計算処理の手順の一例を示すフローチャートである。以下、図20に示す処理をステップ番号に沿って説明する。類似度計算処理は、類似度の比較対象となる2つのデータが入力されたときに実行される。

【0142】

[ステップS101] 入力データ距離算出部221は、変換行列を初期化する。例えば入力データ距離算出部221は、分類対象のデータの「量」以外の変数ごとに、変換行列を生成する。生成される変換行列は、対応する変数に含まれる変数値の数(同一の値の変数値は1つと数える)分の行と列とを有する正方行列である。変換行列の成分には、正規直交条件を満たしていれば、ランダムな値を設定することができる。例えば入力データ距離算出部221は、変換行列内のいくつかの成分の値をランダムに決定し、正規直交条件を満たすように他の成分の値を決定する。

40

【0143】

[ステップS102] 入力データ距離算出部221は、生成した変換行列を用いて、比較対象のデータそれぞれから中間データを生成する。

【0144】

[ステップS103] 入力データ距離算出部221は、中間データ間の類似度を算出する。入力データ距離算出部221は、算出した類似度をメモリに一時的に保存する。

50

【 0 1 4 5 】

[ステップ S 1 0 4] 入力データ距離算出部 2 2 1 は、比較対象のデータの変数を 1 つ選択する。

【 0 1 4 6 】

[ステップ S 1 0 5] 入力データ距離算出部 2 2 1 は、比較対象のデータそれぞれの変数値間の類似度を示す類似度行列を生成する。例えば入力データ距離算出部 2 2 1 は、比較対象のデータそれぞれについて、選択した変数以外の変数を変換行列で変換した中間データを生成する。そして入力データ距離算出部 2 2 1 は、中間データに示される変数値の量の値と、他の中間データに示される変数値の量の値との類似度を、それらの 2 つの変数値間の類似度を示す成分として、類似度行列に設定する。

10

【 0 1 4 7 】

[ステップ S 1 0 6] 入力データ距離算出部 2 2 1 は、類似度行列に基づいて、選択した変数についての新たな変換行列を生成する。

【 0 1 4 8 】

[ステップ S 1 0 7] 入力データ距離算出部 2 2 1 は、すべての変数を選択したか否かを判断する。すべての変数の選択が選択済みとなった場合、処理がステップ S 1 0 8 に進められる。未選択の変数があれば、処理がステップ S 1 0 4 に進められる。

【 0 1 4 9 】

[ステップ S 1 0 8] 入力データ距離算出部 2 2 1 は、各変数について新たに生成した変換行列を用いて、比較対象のデータごとの中間データを生成する。

20

【 0 1 5 0 】

[ステップ S 1 0 9] 入力データ距離算出部 2 2 1 は、ステップ S 1 0 9 で生成した中間データ間の類似度を算出する。

【 0 1 5 1 】

[ステップ S 1 1 0] 入力データ距離算出部 2 2 1 は、処理の終了条件が満たされたか否かを判断する。処理の終了条件とは、例えば類似度が収束したか、またはステップ S 1 0 4 ~ S 1 1 0 のループを所定回数以上繰り返したことである。処理の終了条件が満たされた場合、類似度計算処理が終了する。処理の終了条件が満たされていない場合、入力データ距離算出部 2 2 1 は、変数の選択状態を未選択に初期化して、処理をステップ S 1 0 4 に進める。

30

【 0 1 5 2 】

このような手順で比較対象のデータ間の類似度を計算することができる。以下、図 2 1 ~ 2 4 を参照して、類似度計算の具体例について説明する。

【 0 1 5 3 】

図 2 1 は、初期状態の変換行列を用いた中間データの生成例を示す図である。図 2 1 の例では、第 1 データ 1 0 3 1 と第 2 データ 1 0 3 2 とが比較対象のデータである。まず、第 1 データ 1 0 3 1 の「項 S」の変換行列 1 0 4 1 と「項 R」の変換行列 1 0 4 2 とが初期化され、初期状態の変換行列 1 0 4 1 , 1 0 4 2 を用いて、第 1 データ 1 0 3 1 が中間データ 1 0 5 1 に変換される。同様に、第 2 データ 1 0 3 2 の「項 S」の変換行列 1 0 4 3 と「項 R」の変換行列 1 0 4 4 とが初期化され、初期状態の変換行列 1 0 4 3 , 1 0 4 4 を用いて、第 2 データ 1 0 3 2 が中間データ 1 0 5 2 に変換される。そして、第 1 データ 1 0 3 1 の中間データ 1 0 5 1 と第 2 データ 1 0 3 2 の中間データ 1 0 5 2 との類似度が算出される。図 1 1 の例では、類似度が「0 . 4 0」である。

40

【 0 1 5 4 】

次に、変数「項 S」が選択されたものとする。「項 S」が選択されると、「項 S」の変換行列が更新される。

【 0 1 5 5 】

図 2 2 は、「項 S」の変換行列の更新例を示す図である。第 1 データ 1 0 3 1 について、「項 R」用の変換行列 1 0 4 2 を用いて、「項 S」以外の変数値を変換した中間データ 1 0 5 3 が生成される。同様に第 2 データ 1 0 3 2 について、「項 R」用の変換行列 1 0

50

4 4 を用いて、「項 S」以外の変数値を変換した中間データ 1 0 5 4 が生成される。次に、生成された 2 つの中間データ 1 0 5 3 , 1 0 5 4 それぞれの変数値間の類似度を示す類似度行列 1 0 6 1 が生成される。そして類似度行列 1 0 6 1 に基づいて、掛け合わせることで類似度行列 1 0 6 1 と近似した値を得ることができる 2 つの変換行列 1 0 4 1 a , 1 0 4 3 a が生成される。

【 0 1 5 6 】

次に、変数「項 R」が選択されたものとする。「項 R」が選択されると、「項 R」の変換行列が更新される。

【 0 1 5 7 】

図 2 3 は、「項 R」の変換行列の更新例を示す図である。第 1 データ 1 0 3 1 について、「項 S」用の変換行列 1 0 4 1 a を用いて、「項 R」以外の変数値を変換した中間データ 1 0 5 5 が生成される。同様に第 2 データ 1 0 3 2 について、「項 S」用の変換行列 1 0 4 3 a を用いて、「項 R」以外の変数値を変換した中間データ 1 0 5 6 が生成される。次に、生成された 2 つの中間データ 1 0 5 5 , 1 0 5 6 それぞれの変数値間の類似度を示す類似度行列 1 0 6 2 が生成される。そして類似度行列 1 0 6 2 に基づいて、掛け合わせることで類似度行列 1 0 6 2 と近似した値を得ることができる 2 つの変換行列 1 0 4 2 a , 1 0 4 4 a が生成される。

10

【 0 1 5 8 】

「項 S」と「項 R」とのそれぞれについて、変換行列の更新が終了すると、更新後の変換行列 1 0 4 1 a , 1 0 4 2 a , 1 0 4 3 a , 1 0 4 4 a に基づいて、中間データが生成される。

20

【 0 1 5 9 】

図 2 4 は、更新後の変換行列を用いた中間データの生成例を示す図である。更新後の変換行列 1 0 4 1 a , 1 0 4 2 a を用いて、第 1 データ 1 0 3 1 から中間データ 1 0 5 3 が生成される。同様に、更新後の変換行列 1 0 4 3 a , 1 0 4 4 a を用いて、第 2 データ 1 0 3 2 から中間データ 1 0 5 4 が生成される。そして、生成された中間データ 1 0 5 3 , 1 0 5 4 間の類似度が算出される。図 2 4 の例では、類似度は「0 . 9 1」である。

【 0 1 6 0 】

更新後の変換行列 1 0 4 1 a , 1 0 4 2 a , 1 0 4 3 a , 1 0 4 4 a を用いて生成した中間データ 1 0 5 3 , 1 0 5 4 間の類似度は、更新前の変換行列 1 0 4 1 ~ 1 0 4 4 を用いて生成した中間データ 1 0 5 1 , 1 0 5 2 間の類似度（図 1 8 参照）よりも高くなっている。すなわち変換行列 1 0 4 1 ~ 1 0 4 4 を更新したことにより、生成される中間データ間の類似度が高まっている。更新後の変換行列 1 0 4 1 a , 1 0 4 2 a , 1 0 4 3 a , 1 0 4 4 a に対して、さらに更新処理を行えば、中間データ間の類似度をさらに高めることができる。ただし更新処理を何度も繰り返すと、中間データ間の類似度の上昇度合いが鈍化し、ある程度の類似度に収束する。

30

【 0 1 6 1 】

入力データ距離算出部 2 2 1 は、例えば、類似度の上昇が所定値以下になったとき、最後に算出した中間データ間の類似度を、比較対象の第 1 データ 1 0 3 1 と第 2 データ 1 0 3 2 との類似度に決定する。入力データ距離算出部 2 2 1 は、決定した類似度から距離を算出し、入力データ距離記憶部 2 2 2 に格納する。

40

【 0 1 6 2 】

このように、変換行列を用いて生成した中間データ間の類似度により、複数のデータ間の類似度を決定するようにしたことで、類似度の判定精度が向上し、データ間の距離についての精度も向上する。

【 0 1 6 3 】

以下に、データに含まれる変数値の並べ替えのみによる類似度の計算例である比較例の図 2 5 および図 2 6 と、変換行列を用いて生成した中間データ間の類似度の計算例である図 2 7 および図 2 8 とについて説明する。

【 0 1 6 4 】

50

図 2 5 は、類似度計算の比較例を示す第 1 の図である。図 2 5 に示す第 1 データ 1 0 7 1 と第 2 データ 1 0 7 2 との類似度を計算する場合を想定する。人やモノの間の関係のしかたを分類するとき、第 1 データ 1 0 7 1 を採取した期間に「S 1」の装置が担っていた役割を、第 2 データ 1 0 7 2 を採取した期間では「S 2」の装置が担っている可能性がある。そこで図 2 5 の例では、通信元ホストや通信先ホストを、別の変数値に対応付けて、各データ内のレコードの並べ替えを行っている。

【 0 1 6 5 】

並べ替えにより、変換データ 1 0 7 3 , 1 0 7 4 が生成される。2 つの変換データ 1 0 7 3 , 1 0 7 4 は、通信元ホスト、通信先ホスト、ポートの関係を示す変数値の組み合わせの順番が統一されている。図 2 5 の例では、変換データ 1 0 7 3 , 1 0 7 4 の最上位には、「S ' 1」、「R ' 1」、「P ' 1」の組み合わせを示すレコードが登録され、その次に「S ' 1」、「R ' 1」、「P ' 2」の組み合わせを示すレコードが登録されている。

10

【 0 1 6 6 】

このように変換データ 1 0 7 3 , 1 0 7 4 内に所定の順番で並べられた各レコードの量の値を比較することで、変換データ 1 0 7 3 , 1 0 7 4 間の類似度を算出できる。例えば、量の値を成分とするレベクトル間の内積が、類似度とされる。この場合、第 1 データ 1 0 7 1 と第 2 データ 1 0 7 2 との各変数値に、変換データ 1 0 7 3 , 1 0 7 4 のどの変数値を対応付けるかにより、類似度が変わってくる。そのため、対応付けのすべてのパターンについて変換データ 1 0 7 3 , 1 0 7 4 を生成し、類似度の最大化が図られる。そして、変換データ 1 0 7 3 , 1 0 7 4 から得られる類似度の最大値が、第 1 データ 1 0 7 1 と第 2 データ 1 0 7 2 との類似度と判定される。

20

【 0 1 6 7 】

図 2 6 は、類似度計算の比較例を示す第 2 の図である。図 2 5 に示した方法で類似度を計算したときの第 1 データ 1 0 7 1 と第 2 データ 1 0 7 2 との類似度が「0 . 8 9」である。同じ方法で、第 1 データ 1 0 7 1 と第 3 データ 1 0 7 5 との類似度を計算すると、同じく「0 . 8 9」となる。

【 0 1 6 8 】

ここで、第 1 データ 1 0 7 1 と第 2 データ 1 0 7 2 における通信先ホストとポートとの関係を見ると、「量」の値が「1」以上のレコードに、{ R 1 , P 1 } または { R 2 , P 2 } の組み合わせしか含まれていないことが分かる。それに対して第 3 データ 1 0 7 5 では、「量」の値が「1」以上のレコードのなかに、{ R 1 , P 2 } の組み合わせを含むものがある。そうすると、第 1 データ 1 0 7 1 と第 2 データ 1 0 7 2 とは、通信元ホストが分離・併合された前後での通信ログというだけで、類似の事象に関する通信ログである可能性が高い。それに対して、第 3 データ 1 0 7 5 は、別の事象に関する通信ログであると考えられる。

30

【 0 1 6 9 】

しかし、図 2 5 に示した方法で類似度を計算すると、図 2 6 に示すように、第 1 データ 1 0 7 1 から見たとき、第 2 データ 1 0 7 2 と第 3 データ 1 0 7 5 とのいずれの間も類似度が同じとなる。すなわち、正しく類似度が計算されていない。

40

【 0 1 7 0 】

次に第 1 データ 1 0 7 1 と第 2 データ 1 0 7 2 との類似度、および第 1 データ 1 0 7 1 と第 3 データ 1 0 7 5 との類似度を、第 2 の実施形態に係る方法で計算した場合について、図 2 7 , 図 2 8 を参照して説明する。

【 0 1 7 1 】

図 2 7 は、入力データ距離算出部 2 2 1 による類似度計算例を示す第 1 の図である。図 2 7 には、第 1 データ 1 0 7 1 と第 2 データ 1 0 7 2 との類似度の計算例を示している。第 1 データ 1 0 7 1 について、通信元ホスト、通信先ホスト、ポートそれぞれに対応する変換行列 1 0 8 1 - 2 ~ 1 0 8 3 - 2 が生成されている。これらの変換行列 1 0 8 1 - 2 ~ 1 0 8 3 - 2 を用いて、第 1 データ 1 0 7 1 が中間データ 1 0 9 1 - 2 に変換されてい

50

る。また第2データ1072について、通信元ホスト、通信先ホスト、ポートそれぞれに対応する変換行列1084～1086が生成されている。これらの変換行列1084～1086を用いて、第2データ1072が中間データ1092に変換されている。第1データ1071の中間データ1091と第2データ1072の中間データ1092との類似度は、「0.97」である。

【0172】

図28は、入力データ距離算出部221による類似度計算例を示す第2の図である。図28には、第1データ1071と第3データ1075との類似度の計算例を示している。第1データ1071について、通信元ホスト、通信先ホスト、ポートそれぞれに対応する変換行列1081-3～1083-3が生成されている。これらの変換行列1081-3～1083-3を用いて、第1データ1071が中間データ1091-3に変換されている。第3データ1075について、通信元ホスト、通信先ホスト、ポートそれぞれに対応する変換行列1087～1089が生成されている。これらの変換行列1087～1089を用いて、第3データ1075が中間データ1093に変換されている。第1データ1071の中間データ1091-3と第3データ1075の中間データ1093との類似度は、「0.94」である。

10

【0173】

図27と図28の類似度の計算結果から、第1データ1071は、第3データ1075よりも第2データ1072に類似していることが分かる。すなわち、変数値間の関係を正しく反映させた類似度計算により、類似度の計算制度が向上している。

20

【0174】

しかも図25、図26に示したような方法で類似度の精度を上げようとする、対応付けのすべてのパターンについて類似度を計算することとなり、計算量が膨大となる。

【0175】

例えば、3項目の人またはものがあり、各項目の種類数がそれぞれ「A, B, C」(A, B, Cは1以上の整数)であるものとする。このとき、図25、図26に示した方法で類似度を計算すると、「A!B!C!」の数の組み合わせパターンについて類似度計算を行うこととなる。それに対して、入力データ距離算出部221による手法では、 $(A^2 + B^2 + C^2)ABC$ に比例する計算量となる。これは「A, B, C」がそれぞれ「10, 10, 10」なら、約160,000,000,000,000倍高速となることを意味する。

30

【0176】

以上のように、入力データ距離算出部221による類似度計算方法によれば、類似する事象がログ生成の過程で異なる状態で記録された場合でも、重みづけによる変換により、精度の高い類似度の判定を、効率的に実行することができ、精度の高い距離の判定を行うことができる。

【0177】

〔その他の実施の形態〕

第2の実施形態では、単位期間ごとの通信ログ間の類似度および距離を計算する例を示したが、同じ技術により、他の様々な情報の類似度および距離を計算可能である。

40

【0178】

以上、実施の形態を例示したが、実施の形態で示した各部の構成は同様の機能を有する他のものに置換することができる。また、他の任意の構成物や工程が付加されてもよい。さらに、前述した実施の形態のうちの任意の2以上の構成(特徴)を組み合わせたものであってもよい。

【符号の説明】

【0179】

- 10 分類装置
- 110 収集部
- 112 入力データ記憶部

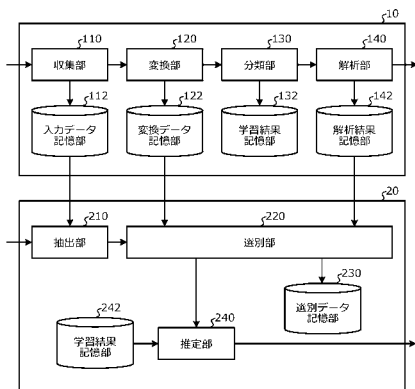
50

- 1 2 0 変換部
- 1 2 2 変換データ記憶部
- 1 3 0 分類部
- 1 3 2 学習結果記憶部
- 1 4 0 解析部
- 1 4 2 解析結果記憶部
- 2 0 要因推定装置
- 2 1 入力データの特徴空間
- 2 2 変換データの特徴空間
- 2 4 近似識別線
- 2 5 - 1 ~ 2 5 - 4 説明ベクトル
- 2 6 点 u における識別要因
- 2 1 0 抽出部
- 2 2 0 選別部
- 2 2 1 入力データ距離算出部
- 2 2 2 入力データ距離記憶部
- 2 2 3 変換データ距離算出部
- 2 2 4 変換データ距離記憶部
- 2 2 5 対象判定部
- 2 3 0 選別データ記憶部
- 2 4 0 推定部
- 2 4 2 学習結果記憶部

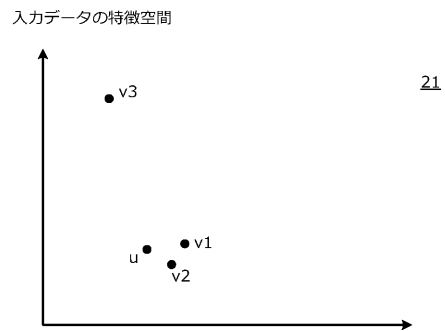
10

20

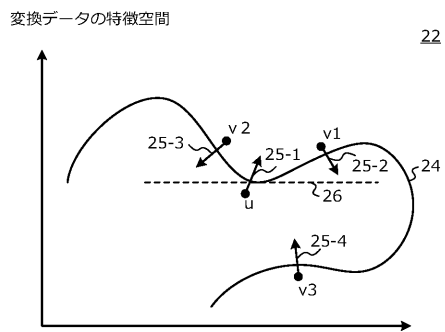
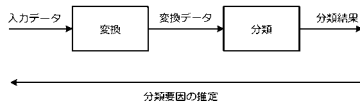
【 図 1 】



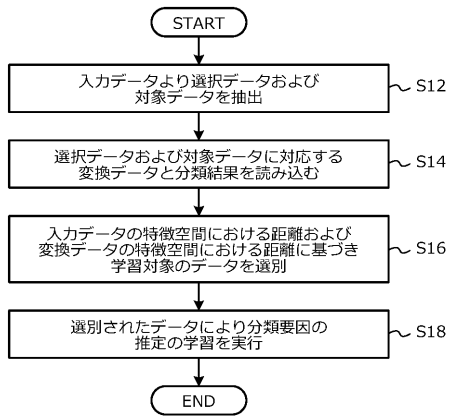
【 図 3 】



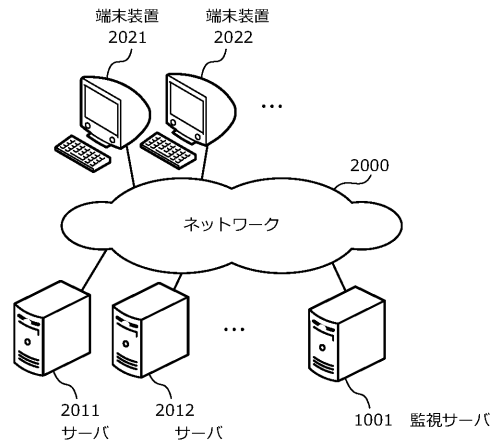
【 図 2 】



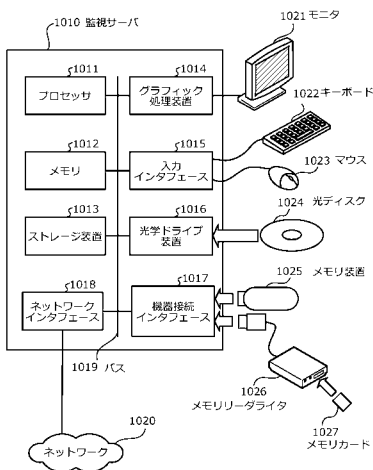
【 図 4 】



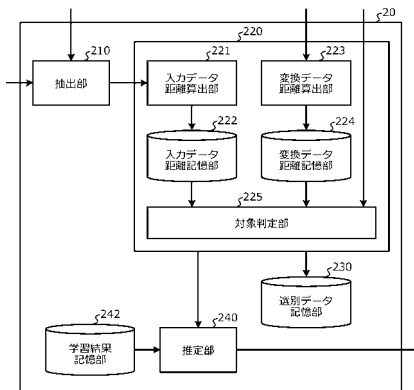
【 図 5 】



【 図 6 】



【 図 7 】



【 図 8 A 】

801 入力データ				803 順序付け				805 中間データ			
通信元ホスト	通信先ホスト	ポート	重			通信元ホスト	通信先ホスト	ポート	重		
S3	R1	P1	1	S3→S1	S1	R'1	P'1	1			
S3	R2	P1	1	S1→S2	S1	R'2	P'1	1			
S3	R3	P1	1	S2→S3	S1	R'3	P'1	1			
S1	R3	P3	2	R1→R1	S1	R'3	P'2	2			
S2	R1	P1	1	R2→R2	S2	R'3	P'1	1			
				R3→R3	S3	R'1	P'1	1			
				P1→P1							
				P3→P2							

【図 8 B】

図 8B は、入力データ、中間データ、および変換データのテーブルを示しています。各テーブルは送信元ホスト、送信先ホスト、ポート、および単位 (単位) を示しています。

811 入力データ

送信元ホスト	送信先ホスト	ポート	単位
S3	R1	P1	1
S3	R2	P1	1
S2	R3	P1	1
S2	R3	P3	1
S2	R2	P1	1

815 中間データ

送信元ホスト	送信先ホスト	ポート	単位
S1	R1	P1	1
S1	R2	P1	1
S2	R3	P1	1
S2	R3	P2	1
S2	R2	P1	1

821 入力データ

送信元ホスト	送信先ホスト	ポート	単位
S2	R1	P1	1
S2	R2	P1	1
S1	R3	P1	1
S2	R3	P2	1
S1	R1	P1	1
S2	R1	P2	1

825 中間データ

送信元ホスト	送信先ホスト	ポート	単位
S1	R1	P1	1
S1	R2	P1	1
S1	R3	P1	1
S2	R1	P1	1
S2	R1	P2	1

831 入力データ

送信元ホスト	送信先ホスト	ポート	単位
S3	R1	P1	1
S3	R2	P1	1
S2	R3	P1	1
S2	R3	P3	1
S2	R2	P1	1

835 中間データ

送信元ホスト	送信先ホスト	ポート	単位
S1	R1	P1	1
S1	R2	P1	1
S2	R3	P1	1
S2	R3	P2	1
S2	R1	P1	1

【図 9 B】

図 9B は、入力データ、中間データ、および変換データのテーブルを示しています。各テーブルは送信元ホスト、送信先ホスト、ポート、および単位 (単位) を示しています。

911 入力データ

送信元ホスト	送信先ホスト	ポート	単位
S3	R1	P1	1
S3	R2	P1	1
S2	R3	P1	1
S2	R3	P2	1
S2	R2	P1	1

915 変換データ

送信元ホスト	送信先ホスト	ポート	単位
S2	R1	P1	1
S2	R2	P1	1
S1	R3	P1	1
S2	R3	P2	1
S2	R2	P1	1

921 入力データ

送信元ホスト	送信先ホスト	ポート	単位
S2	R1	P1	1
S2	R2	P1	1
S1	R3	P1	1
S2	R3	P2	1
S1	R1	P1	1
S2	R1	P2	1

925 変換データ

送信元ホスト	送信先ホスト	ポート	単位
S2	R1	P1	1
S2	R2	P1	1
S1	R3	P1	1
S2	R3	P2	1
S2	R1	P1	1
S2	R1	P2	1

931 入力データ

送信元ホスト	送信先ホスト	ポート	単位
S3	R1	P1	1
S3	R2	P1	1
S2	R3	P1	1
S2	R3	P2	1
S2	R2	P1	1

935 変換データ

送信元ホスト	送信先ホスト	ポート	単位
S2	R1	P1	1
S2	R2	P1	1
S1	R3	P1	1
S2	R3	P2	1
S2	R1	P1	1
S2	R1	P2	1

【図 10】

	v1	v2	v3
入力データ特徴空間	5	4	9
変換データ特徴空間	9	8	9
合計	14	12	18

【図 9 A】

図 9A は、入力データ、中間データ、および変換データのテーブルを示しています。各テーブルは送信元ホスト、送信先ホスト、ポート、および単位 (単位) を示しています。

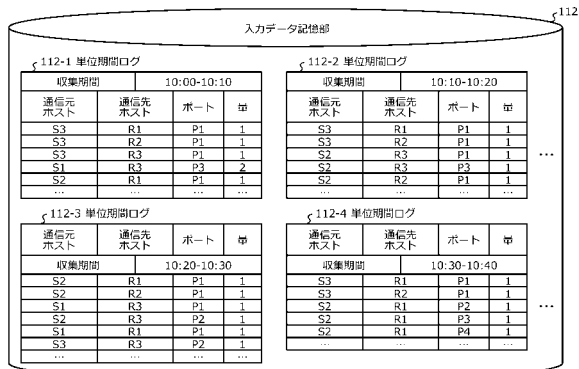
901 入力データ

送信元ホスト	送信先ホスト	ポート	単位
S3	R1	P1	1
S3	R2	P1	1
S3	R3	P1	1
S1	R3	P3	2
S2	R1	P1	1

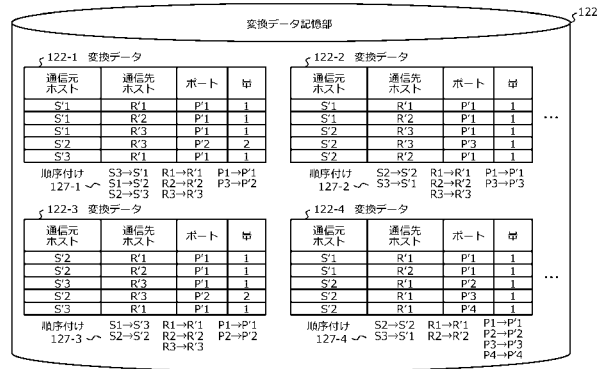
905 変換データ

送信元ホスト	送信先ホスト	ポート	単位
S'1	R'1	P'1	1
S'1	R'2	P'1	1
S'1	R'3	P'1	1
S'2	R'3	P'2	2
S'3	R'1	P'1	1

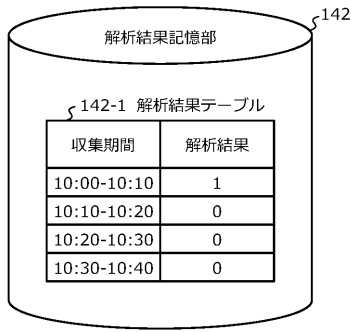
【図 11】



【図 12】



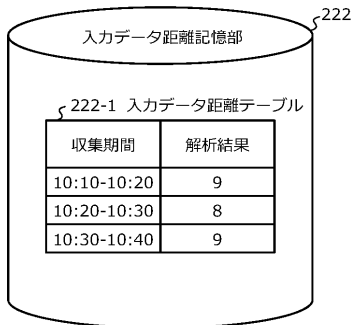
【 図 1 3 】



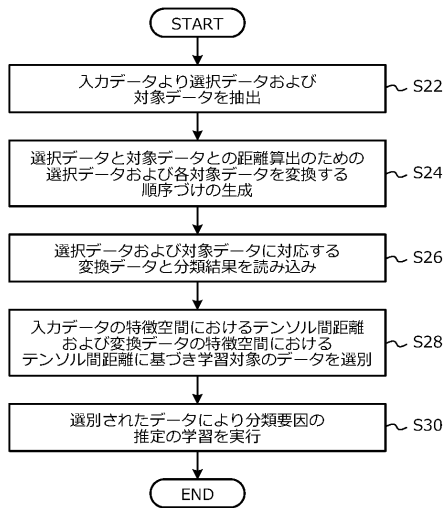
【 図 1 5 】



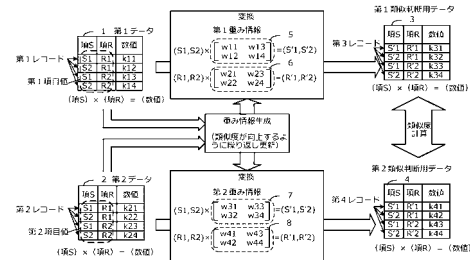
【 図 1 4 】



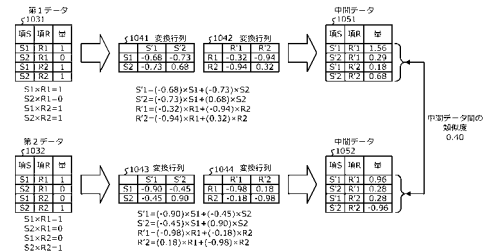
【 図 1 6 】



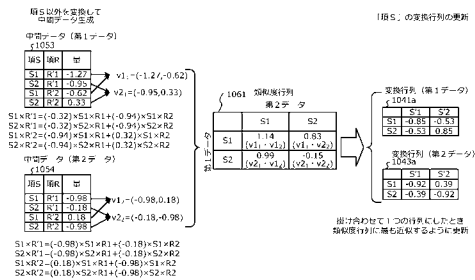
【 図 1 7 】



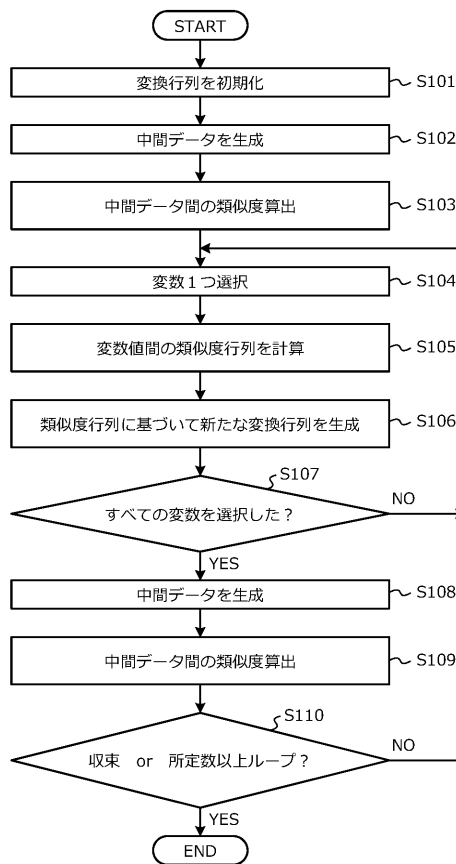
【 図 1 8 】



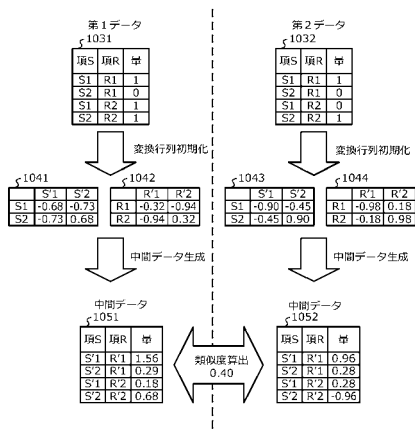
【図 19】



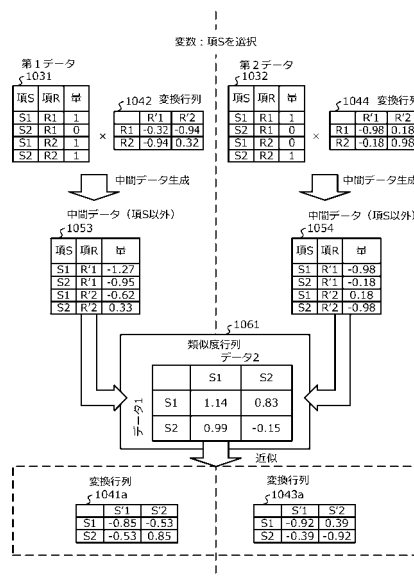
【図 20】



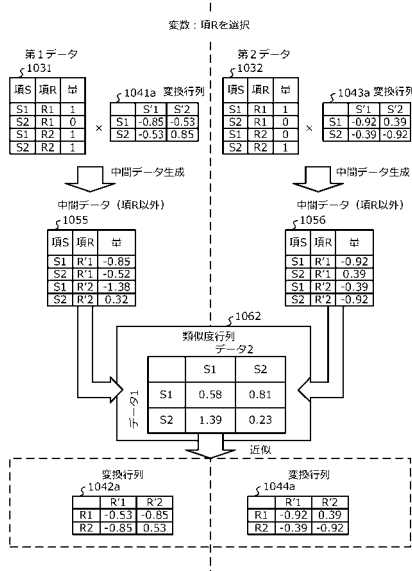
【図 21】



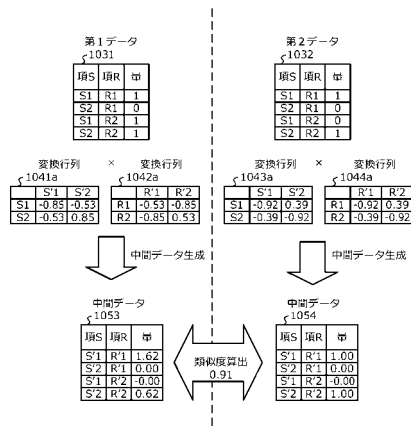
【図 22】



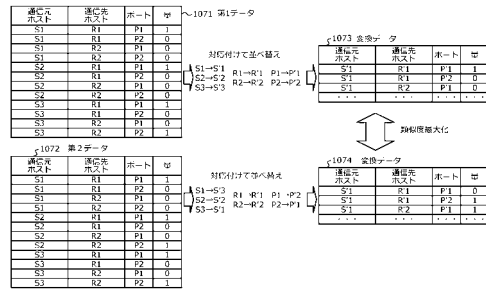
【図 2 3】



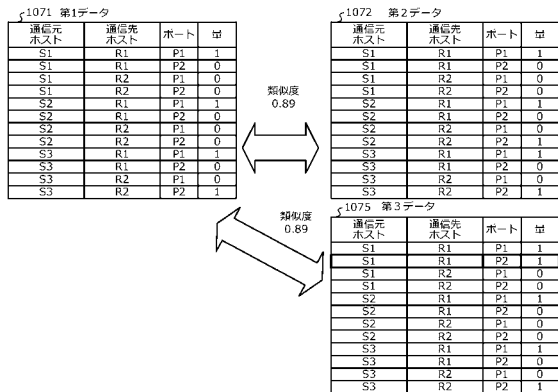
【図 2 4】



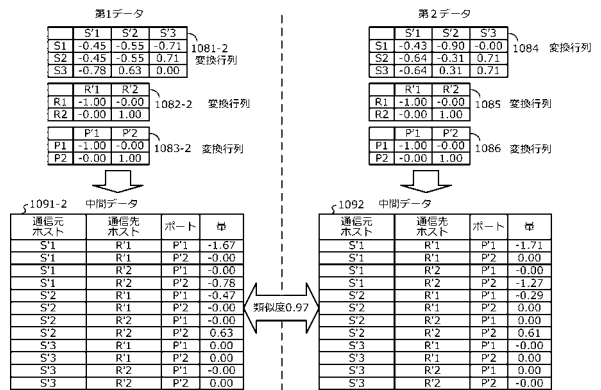
【図 2 5】



【図 2 6】



【図 2 7】



【 図 2 8 】

