

(19) 日本国特許庁(JP)

(12) 特 許 公 報(B2)

(11) 特許番号

特許第5391637号  
(P5391637)

(45) 発行日 平成26年1月15日(2014.1.15)

(24) 登録日 平成25年10月25日(2013.10.25)

(51) Int.Cl. F I  
**G06F 17/30 (2006.01)** G O 6 F 17/30 2 1 O D  
 G O 6 F 17/30 3 5 O C

請求項の数 24 (全 27 頁)

(21) 出願番号 特願2008-264442 (P2008-264442)  
 (22) 出願日 平成20年10月10日(2008.10.10)  
 (65) 公開番号 特開2010-92432 (P2010-92432A)  
 (43) 公開日 平成22年4月22日(2010.4.22)  
 審査請求日 平成23年9月8日(2011.9.8)

(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  
 (74) 代理人 100103090  
 弁理士 岩壁 冬樹  
 (74) 代理人 100124501  
 弁理士 塩川 誠人  
 (72) 発明者 黒岩 由希子  
 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 野崎 大進

最終頁に続く

(54) 【発明の名称】 データ類似度計算システム、データ類似度計算方法およびデータ類似度計算プログラム

(57) 【特許請求の範囲】

【請求項1】

データの特徴を示す属性とデータの類別を示すクラスとを含むデータの集合から、データの重複を許して前記集合の部分集合を複数回生成する部分集合生成手段と、

部分集合が生成される毎に、属性からクラスを判定するルールである分類器を、前記部分集合に基づいて生成する分類器生成手段と、

分類器が生成される毎に、分類器を用いて、前記部分集合に属する個々のデータのクラスを判定するクラス判定手段と、

データの集合の部分集合が生成され、クラス判定手段が当該部分集合に属する個々のデータのクラスを判定したときに、同一のクラスと判定されたデータ同士の類似度に値を加算する類似度算出手段とを備える

ことを特徴とするデータ類似度計算システム。

【請求項2】

データ同士の類似度に基づいて、データの集合に属するデータを複数のグループに分類するデータグループ化手段を備える

請求項1に記載のデータ類似度計算システム。

【請求項3】

データグループ化手段は、データの集合に属する個々のデータをそれぞれ別々のグループに分類し、互いに異なる二つのグループに属するデータ同士の類似度を求め、前記類似度が最大となる二つのグループを併合することを繰り返し、グループの総数を目標数まで

減少させる

請求項 2 に記載のデータ類似度計算システム。

【請求項 4】

データ集合に属するデータを、特定の属性の属性値に応じて、グループに分類する属性データ分類手段と、

データグループ化手段によって分類されたデータのグループと、属性データ分類手段によって分類されたデータのグループとの関係に基づいて、類似度算出に対する前記特定の属性の関連度を計算する関連度計算手段とを備える

請求項 2 または請求項 3 に記載のデータ類似度計算システム。

【請求項 5】

部分集合生成手段は、データの集合からデータをランダムサンプリングすることによって、前記集合の部分集合を生成する

請求項 1 から請求項 4 のうちのいずれか 1 項に記載のデータ類似度計算システム。

【請求項 6】

類似度算出手段は、クラス判定手段によって特定のクラスと判定されたデータ同士の類似度に対してのみ値を加算する

請求項 1 から請求項 5 のうちのいずれか 1 項に記載のデータ類似度計算システム。

【請求項 7】

分類器生成手段は、部分集合に属するデータのうち所与のクラスが特定のクラスであるデータを加重して分類器を生成する

請求項 1 から請求項 6 のうちのいずれか 1 項に記載のデータ類似度計算システム。

【請求項 8】

部分集合生成手段は、少なくとも顧客の特徴または販売条件を属性とし顧客の行動をクラスとするデータの集合から、当該集合の部分集合を生成する

請求項 1 から請求項 7 のうちのいずれか 1 項に記載のデータ類似度計算システム。

【請求項 9】

部分集合生成手段が、データの特徴を示す属性とデータの類別を示すクラスとを含むデータの集合から、データの重複を許して前記集合の部分集合を複数回生成する部分集合生成ステップと、

分類器生成手段が、部分集合が生成される毎に、属性からクラスを判定するルールである分類器を、前記部分集合に基づいて生成する分類器生成ステップと、

クラス判定手段が、分類器が生成される毎に、分類器を用いて、前記部分集合に属する個々のデータのクラスを判定するクラス判定ステップと、

類似度算出手段が、データの集合の部分集合が生成され、クラス判定ステップで当該部分集合に属する個々のデータのクラスを判定したときに、同一のクラスと判定されたデータ同士の類似度に値を加算する類似度算出ステップとを含む

ことを特徴とするデータ類似度計算方法。

【請求項 10】

データグループ化手段が、データ同士の類似度に基づいて、データの集合に属するデータを複数のグループに分類するデータグループ化ステップを含む

請求項 9 に記載のデータ類似度計算方法。

【請求項 11】

データグループ化手段が、データグループ化ステップで、データの集合に属する個々のデータをそれぞれ別々のグループに分類し、互いに異なる二つのグループに属するデータ同士の類似度を求め、前記類似度が最大となる二つのグループを併合することを繰り返し、グループの総数を目標数まで減少させる

請求項 10 に記載のデータ類似度計算方法。

【請求項 12】

属性データ分類手段が、データ集合に属するデータを、特定の属性の属性値に応じて、グループに分類する属性データ分類ステップと、

	10
	20
	30
	40
	50

関連度計算手段が、データグループ化ステップで分類されたデータのグループと、属性データ分類ステップで分類されたデータのグループとの関係に基づいて、類似度算出に対する前記特定の属性の関連度を計算する関連度計算ステップとを備える

請求項 10 または 請求項 11 に記載のデータ類似度計算方法。

【請求項 13】

部分集合生成手段が、部分集合生成ステップで、データの集合からデータをランダムサンプリングすることによって、前記集合の部分集合を生成する

請求項 9 から請求項 12 のうちのいずれか 1 項に記載のデータ類似度計算方法。

【請求項 14】

類似度算出手段が、類似度算出ステップで、クラス判定ステップで特定のクラスと判定されたデータ同士の類似度に対してのみ値を加算する

請求項 9 から請求項 13 のうちのいずれか 1 項に記載のデータ類似度計算方法。

【請求項 15】

分類器生成手段が、分類器生成ステップで、部分集合に属するデータのうち所与のクラスが特定のクラスであるデータを加重して分類器を生成する

請求項 9 から請求項 14 のうちのいずれか 1 項に記載のデータ類似度計算方法。

【請求項 16】

部分集合生成手段が、部分集合生成ステップで、少なくとも顧客の特徴または販売条件を属性とし顧客の行動をクラスとするデータの集合から、当該集合の部分集合を生成する

請求項 9 から請求項 15 のうちのいずれか 1 項に記載のデータ類似度計算方法。

【請求項 17】

コンピュータに、

データの特徴を示す属性とデータの類別を示すクラスとを含むデータの集合から、データの重複を許して前記集合の部分集合を複数回生成する部分集合生成処理、

部分集合が生成される毎に、属性からクラスを判定するルールである分類器を、前記部分集合に基づいて生成する分類器生成処理、

分類器が生成される毎に、分類器を用いて、前記部分集合に属する個々のデータのクラスを判定するクラス判定処理、および、

データの集合の部分集合が生成され、クラス判定処理で当該部分集合に属する個々のデータのクラスを判定したときに、同一のクラスと判定されたデータ同士の類似度に値を加算する類似度算出処理

を実行させるためのデータ類似度計算プログラム。

【請求項 18】

コンピュータに、

データ同士の類似度に基づいて、データの集合に属するデータを複数のグループに分類するデータグループ化処理

を実行させる請求項 17 に記載のデータ類似度計算プログラム。

【請求項 19】

コンピュータに、

データグループ化処理で、データの集合に属する個々のデータをそれぞれ別々のグループに分類させ、互いに異なる二つのグループに属するデータ同士の類似度を求めさせ、前記類似度が最大となる二つのグループを併合することを繰り返させ、グループの総数を目標数まで減少させる

請求項 18 に記載のデータ類似度計算プログラム。

【請求項 20】

コンピュータに、

データ集合に属するデータを、特定の属性の属性値に応じて、グループに分類する属性データ分類処理、および、

データグループ化処理で分類されたデータのグループと、属性データ分類処理で分類されたデータのグループとの関係に基づいて、類似度算出に対する前記特定の属性の関連度

10

20

30

40

50

を計算する関連度計算処理

を実行させる請求項 18 または請求項 19 に記載のデータ類似度計算プログラム。

【請求項 21】

コンピュータに、

部分集合生成処理で、データの集合からデータをランダムサンプリングすることによって、前記集合の部分集合を生成させる

請求項 17 から請求項 20 のうちのいずれか 1 項に記載のデータ類似度計算プログラム

。

【請求項 22】

コンピュータに、

類似度算出処理で、クラス判定処理で特定のクラスと判定されたデータ同士の類似度に対してのみ値を加算させる

請求項 17 から請求項 21 のうちのいずれか 1 項に記載のデータ類似度計算プログラム

。

【請求項 23】

コンピュータに、

分類器生成処理で、部分集合に属するデータのうち所与のクラスが特定のクラスであるデータを加重して分類器を生成させる

請求項 17 から請求項 22 のうちのいずれか 1 項に記載のデータ類似度計算プログラム

。

【請求項 24】

コンピュータに、

部分集合生成処理で、少なくとも顧客の特徴または販売条件を属性とし顧客の行動をクラスとするデータの集合から、当該集合の部分集合を生成させる

請求項 17 から請求項 23 のうちのいずれか 1 項に記載のデータ類似度計算プログラム

。

【発明の詳細な説明】

【技術分野】

【0001】

本発明は、データ間の類似度を計算するデータ類似度計算システム、データ類似度計算方法およびデータ類似度計算プログラムに関する。

【背景技術】

【0002】

小売店では、POS (Point of Sales) システムが用いられており、収集された商品の売上データは様々な活用されている。そして、顧客の特徴に基づいて顧客を分類して、売上データを用いて種々の分析を行うシステムが提案されている。例えば、特許文献 1 には、分析者の意志に従って顧客を分類し、顧客の購買履歴状況を確認可能とするシステムが記載されている。

【0003】

また、特許文献 2 には、POS データを用いて消費者の併売傾向を把握するシステムが記載されている。

【0004】

また、特許文献 3 には、各種カテゴリにおける複数の商品を関連づけて登録するための複数の仮想 BOX を有するデータベースを備え、顧客に対する推奨商品グループを提示できるシステムが記載されている。

【0005】

また、データの類似度を求める装置が特許文献 4 に記載されている。特許文献 4 に記載された装置は、予め保持されている事例と、与えられた事例 (クエリ) との類似度である事例間類似度を計算する。また、保持されている事例 (例えば、ある地点の気象データ) にはクラス (例えば、気象データ計測時点から 3 時間後の天気) が定められている。特許

10

20

30

40

50

文献4に記載された装置は、保持されている事例から相関ルールを導出する。相関ルールは、条件（事例）に応じた結論（クラス）を導くルールである。そして、相関ルールを用いて、事例間類似度から総合類似度を計算する。

【0006】

特許文献4には、「 $A_1, A_2, \dots, A_k, B$ 」という形式で表される相関ルールが記載されている。そして、「 $A_1, A_2, \dots, A_k, B$ 」はアイテム集合と称されている。また、アイテム集合 $A_1, A_2, \dots, A_k, B$ を含むレコード数の全レコードに対する比率を、相関ルールの支持度と呼び、アイテム集合 $A_1, A_2, \dots, A_k$ を含むレコードの中で、アイテム $B$ を含むレコードの割合を相関ルールの確信度と呼ぶ。特許文献4に記載された装置は、支持度および相関度が予め定められた各々の下限値以上となるような相関ルールを抽出する。

10

【0007】

また、自動的にデータを分類する装置が特許文献5に記載されている。特許文献5に記載されたデータ分類装置は、注目パターンデータとの類似度の高い順に複数のパターンデータを取りだし、近傍クラスタを決定する。

【0008】

また、ルールを用いて予測を行う装置として、特許文献6に記載された装置がある。特許文献6に記載された装置は、ある地域（第1の地域）に生じた事象から論理規則を発生し、その規則を他の地域（第2の地域）に適用して、第2の地域でその事象が生じる傾向を予測する。

20

【0009】

【特許文献1】特開平9-101984号公報（段落0014-0050）

【特許文献2】特開2007-94592号公報（段落0024）

【特許文献3】特開2003-337886号公報（段落0010, 0015, 0016）

【特許文献4】特開2002-149697号公報（段落0018-0065）

【特許文献5】特開2003-256443号公報（段落0090-0093）

【特許文献6】特開2004-126757号公報（段落0033, 0051）

【発明の開示】

【発明が解決しようとする課題】

30

【0010】

特許文献1に記載されたシステムでは、システムのユーザである分析者が、手動で顧客を分類する必要があった。そのため、顧客の購買行動から顧客の分類を自動的に行うことができなかった。また、特許文献3に記載された装置においても、各種カテゴリにおける複数の商品の関連付けを、顧客自身や店員等の商品観察者が行う必要があった。

【0011】

データの分類を行うために、データ間の類似度を計算し、その類似度を用いてデータを分類することが考えられる。特許文献4に記載された装置は、相関ルールを導出し、総合類似度を計算する。しかし、特許文献4に記載された装置のように相関ルールを導出すると、関連のある属性（特許文献4におけるアイテム）がデータ中に多く存在する場合に、それらの関連のある属性の影響を強く受け、適切な類似度計算を行えない。

40

【0012】

例えば、図12に示すデータ番号1~10の10個のデータが与えられているとする。図12に示すデータにおいて、データ番号1~4およびデータ番号6は、属性 $b, c, d, e$ の値が共通であり、データの類別を表すクラスも共通である。従って、データの属性値を条件部としてクラスを結論部とする相関ルールを抽出する場合、多くのデータ間で属性値の組み合わせが共通となっている属性 $b, c, d, e$ の影響を強く受け、この結果、類似度も、関連のある特定の属性に影響されてしまう。すなわち、ある属性の属性値がある値に合致しているデータ間の類似度が比較的高くなり、合致していないデータ間の類似度は比較的低くなる傾向が出てしまう。

50

## 【0013】

そこで、本発明は、データを分類するためにデータ間の類似度を求める際に、データ中に互いに関連のある属性が多く存在してもそれらの関連のある属性の影響を強く受けることなく類似度を求めることができるデータ類似度計算システム、データ類似度計算方法およびデータ類似度計算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】

## 【0014】

本発明によるデータ類似度計算システムは、データの特徴を示す属性とデータの類別を示すクラスとを含むデータの集合から、データの重複を許してその集合の部分集合を複数回生成する部分集合生成手段と、部分集合が生成される毎に、属性からクラスを判定するルールである分類器を、部分集合に基づいて生成する分類器生成手段と、分類器が生成される毎に、分類器を用いて、部分集合に属する個々のデータのクラスを判定するクラス判定手段と、データの集合の部分集合が生成され、クラス判定手段がその部分集合に属する個々のデータのクラスを判定したときに、同一のクラスと判定されたデータ同士の類似度に値を加算する類似度算出手段とを備えることを特徴とする。

10

## 【0015】

本発明によるデータ類似度計算方法は、部分集合生成手段が、データの特徴を示す属性とデータの類別を示すクラスとを含むデータの集合から、データの重複を許してその集合の部分集合を複数回生成する部分集合生成ステップと、分類器生成手段が、部分集合が生成される毎に、属性からクラスを判定するルールである分類器を、部分集合に基づいて生成する分類器生成ステップと、クラス判定手段が、分類器が生成される毎に、分類器を用いて、部分集合に属する個々のデータのクラスを判定するクラス判定ステップと、類似度算出手段が、データの集合の部分集合が生成され、クラス判定ステップでその部分集合に属する個々のデータのクラスを判定したときに、同一のクラスと判定されたデータ同士の類似度に値を加算する類似度算出ステップとを含むことを特徴とする。

20

## 【0016】

本発明によるデータ類似度計算プログラムは、コンピュータに、データの特徴を示す属性とデータの類別を示すクラスとを含むデータの集合から、データの重複を許してその集合の部分集合を複数回生成する部分集合生成処理、部分集合が生成される毎に、属性からクラスを判定するルールである分類器を、部分集合に基づいて生成する分類器生成処理、分類器が生成される毎に、分類器を用いて、部分集合に属する個々のデータのクラスを判定するクラス判定処理、および、データの集合の部分集合が生成され、クラス判定処理でその部分集合に属する個々のデータのクラスを判定したときに、同一のクラスと判定されたデータ同士の類似度に値を加算する類似度算出処理を実行させることを特徴とする。

30

## 【発明の効果】

## 【0017】

本発明によれば、データを分類するためにデータ間の類似度を求める際に、データ中に互いに関連のある属性が多く存在してもそれらの関連のある属性の影響を強く受けることなく類似度を求めることができる。

## 【発明を実施するための最良の形態】

40

## 【0018】

以下、本発明の実施形態を図面を参照して説明する。

ここでは、商品やサービスの顧客に関するデータ間の類似度を求める場合を例にして本発明の実施形態を説明する。なお、有料で商品やサービスを利用する者だけでなく、無料で商品やサービスを利用する者や、今後商品やサービスを利用すると思われる人も顧客と呼ぶ。無料の場合は、売価は0とする。また、商品またはサービスの概念としては、製品等の完成品だけでなく、製品の機能を維持するための最小単位も含まれる。さらに、ある商品またはサービスのカテゴリを1つの商品と考えて適用することもできる。

## 【0019】

実施形態1.

50

図1は、本発明の第1の実施形態の例を示すブロック図である。本発明によるデータ類似度計算システム1は、部分集合生成部11と、分類器生成部12と、自己評価部13と、類似度算出部14とを備える。

【0020】

部分集合生成部11は、データの集合から、その集合の部分集合を生成する。データの集合は、例えば、データ類似度計算システム1に設けられたキーボード等の入力装置（図示せず）を介して入力されればよいが、データの集合の入力態様は特に限定されない。以下、データの集合とは、データ全体の集合を意味する。

【0021】

データの集合に含まれる個々のデータは、データの特徴を示す属性と、データの類別を示すクラスとを含む。ここでは、少なくとも顧客の特徴または販売条件をデータの特徴とし、顧客の行動をクラスとする場合を例にして説明する。図2は、データの例を示す説明図である。図2に示す例では、「天気」、「気温」、「湿度」、「風が強いかな否か」を属性とする場合を示している。これらの属性は、販売条件であるが、「年齢」、「性別」等の顧客の特徴を属性としてもよく、顧客の特徴および販売条件がいずれも属性となってもよい。また、図2に示す例では、顧客の行動は、「購入する」と「購入しない」の2種類であるものとし、データも「購入する」および「購入しない」という2種類のクラスに類別される場合を例にする。図2に例示するデータは、属性が示す条件の下で、ある商品またはサービスを顧客が購入したか否かを表している。例えば、図1に示す1番目のデータは、「晴れ」、気温「29度」、湿度「85%」、「風が強くない」という条件の下で、顧客が購入しなかったということを表している。

【0022】

部分集合生成部11は、このようなデータの集合から一部のデータを選択することで、データの集合の部分集合を生成する。また、部分集合生成部11は、部分集合を複数回生成する。このとき、データの重複を許して部分集合を生成する。すなわち、部分集合生成部11によって生成されたある部分集合と別の部分集合に、同じデータが属していてもよい。

【0023】

分類器生成部12は、データの集合の部分集合が生成される毎に、その部分集合に基づいて分類器を生成する。分類器は、属性からクラスを判定するルールである。部分集合と、その部分集合に基づいて生成された分類器とが対応する。

【0024】

自己評価部13は、分類器が生成される毎に、生成された分類器を用いて、その分類器に対応する部分集合に属する個々のデータのクラスを判定する。この判定は、分類器を用いて、個々のデータに含まれる属性からクラスを予測する処理であるということもできる。なお、自己評価部13が判定したクラス（換言すれば、予測したクラス）と、データにおいて予め定められたクラスとが合致するとは限らない。

【0025】

類似度算出部14は、自己評価部13で判定されたクラスに基づいて、部分集合に属する個々のデータ間の類似度を算出する。データの集合に属する任意のデータ同士の組み合わせに対して、予め類似度の初期値0が設定される。類似度算出部14は、部分集合生成部11によって部分集合が生成され、自己評価部13がその部分集合に属する個々のデータのクラスを判定したときに、同一のクラスと判定されたデータ間の類似度に値を加算していくことによって、各データ間の類似度を算出する。

【0026】

部分集合生成部11、分類器生成部12、自己評価部13および類似度算出部14は、例えば、プログラム（データ類似度計算プログラム）に従って動作するCPUによって実現される。その場合、プログラムは、例えばデータ類似度計算システム1が備えるプログラム記憶装置（図示せず）に記憶され、CPUがプログラムを読み込み、そのプログラムに従って部分集合生成部11、分類器生成部12、自己評価部13および類似度算出部1

10

20

30

40

50

4として動作すればよい。

【0027】

次に、動作について説明する。

図3は、第1の実施形態のデータ類似度計算システムの処理経過の例を示すフローチャートである。例えば、データ類似度計算システム1に設けられたキーボード等の入力装置（図示せず）を介して、データの集合が入力されると、データ類似度計算システムは以下のように動作する。ただし、データの集合に属する個々のデータ（個別データと称する）の数をNとする。また、データの集合の部分集合を生成する回数をTとする。Tは、例えば、100、500、1000等の数であるが、Tはこれらの値に限定されない。また、部分集合生成部11が1つの部分集合に属する個別データとしてデータの集合から選択する個別データの数をMとする。Mは、例えば、Nの1%、5%、または10%等の値とすればよいが、Mはこれらの値に限定されない。部分集合の生成を繰り返す繰り返し回数であるTや、部分集合に含める個別データ数であるMは、それぞれ、キーボード等の入力装置（図示せず）を介して、データ類似度計算システムのユーザによって入力されてもよい。あるいは、他の態様でT、Mが指定されてもよい。

10

【0028】

また、データの集合に属する個別データを順番に指定するための第1の変数を*i*とし、第2の変数を*j*とする。*i*は、1 ≤ *i* ≤ Nを満たす整数であり、同様に、*j*は、1 ≤ *j* ≤ Nを満たす整数である。*i*、*j*を指定することでそれぞれ個別データを指定することができ、*i*、*j*の組み合わせによって一対の個別データの組を指定することができる。*i*、*j*によって指定される1対の個別データの類似度をSim(*i*、*j*)と記す。

20

【0029】

最初に、部分集合生成部11は、1 ≤ *i* ≤ Nの範囲の*i*と、1 ≤ *j* ≤ Nの範囲の*j*とによって定められる*i*と*j*の組み合わせを順次定め、その各組み合わせに応じた各類似度Sim(*i*、*j*)を0に初期化する（ステップA1）。すなわち、*i*番目の個別データと*j*番目の個別データとの類似度となるSim(*i*、*j*)を、*i*、*j*の組み合わせ毎に0に初期化する。換言すれば、部分集合生成部11は、データの集合に属する全ての個別データから得られる一対の個別データの各組について、個別データ間の類似度を0に初期化する。

【0030】

次に、部分集合生成部11は、データの集合から部分集合を生成した回数を表す変数*t*を0に初期化する（ステップA2）。

30

【0031】

続いて、部分集合生成部11は、変数*t*と、Tとを比較し、*t*がT未満であるか否かを判定する（ステップA3）。変数*t*がT以上であるならば（ステップA3におけるNo）、処理を終了する。また、変数*t*がT未満であるならば（ステップA3におけるYes）、部分集合生成部11は、データの集合からM個の個別データを選択することによって部分集合を生成する（ステップA4）。すなわち、部分集合の要素となるM個の個別データをデータの集合から選択すればよい。ステップA4の処理は*t*の値がインクリメントされる毎に行われるが、2回目以降の部分集合の生成時において、新たに生成する部分集合に属する個別データと、既に生成された部分集合に属する個別データとが重複していてもよい。このように、部分集合生成部11は、部分集合同士での個別データの重複を許して部分集合を生成する。

40

【0032】

ステップA4において、データの集合から個別データを選択する際に、例えば、個別データをランダムサンプリングしてもよい。すなわち、部分集合生成部11は、データの集合から1個の個別データをランダムに選択することをM回繰り返すことによって、M個の個別データを選択してもよい。

【0033】

ただし、M個の個別データの選択方法はランダムサンプリングに限定されない。例えば

50



、部分集合生成部 1 1 は、最初にステップ A 4 を実行するときには、1 番目から M 番目までの個別データを選択し、次にステップ A 4 を実行するときには、M + 1 番目から 2・M 番目までの個別データを選択するというように、個別データに予め定められた順番に従って、ステップ A 4 毎に M 個ずつ個別データを選択してもよい。

#### 【0034】

なお、部分集合に含める個別データは、データの集合に属する個別データに含まれる属性のうち一部の属性と、クラスを抽出したデータであってもよい。例えば、図 2 に示す例では、個別データには「天気」、「天候」、「湿度」、「風が強いが否か」という 4 つの属性が含まれているが、部分集合に含める個別データを選択するときには、この 4 つの個別データのうちの一部分（例えば、「天気」、「天候」のみ）と、クラスとを抽出し、その属性およびクラスからなるデータを部分集合の要素としてもよい。データの集合に属する個別データに含まれる属性が A 個であるとき、部分集合に属する個別データに含まれる属性をいくつにするか、また、どの属性を部分集合に属する個別データに含めるかについては、例えば、キーボード等の入力装置（図示せず）を介して、データ類似度計算システムのユーザによって入力されてもよいし、他の態様で指定されてもよい。

10

#### 【0035】

ステップ A 4 の後、分類器生成部 1 2 は、直前のステップ A 4 で生成された部分集合に基づいて分類器を生成する（ステップ A 5）。分類器として、例えば決定器やサポートベクタマシンを生成すればよい。どのような分類器を用いるかについては、キーボード等の入力装置（図示せず）を介して、データ類似度計算システムのユーザによって入力されてもよいし、他の態様で指定されてもよい。

20

#### 【0036】

以下、分類器として決定木を用いる場合を例にして、分類器の具体例を説明する。本例では、「購入する、購入しない」をクラスとし、「天気」、「気温」、「湿度」、「風が強いが否か」を属性とする場合を例にして説明する。また、ここでは、図 2 に示す 1 4 個の個別データからなる部分集合が生成されている場合を例にする。図 4 は、分類器の例を示す説明図であり、ここでは、決定木を分類器としている。図 2 に示すような既知の属性およびクラスの組み合わせがあれば、その既知の属性およびクラスから決定木を生成することができる。以下の説明において、属性の値（すなわち属性値）に応じて、決定木におけるノードを枝分かれさせることを分割という。

30

#### 【0037】

図 2 に示す例では、各行が、「購入する」または「購入しない」というクラスが付された個別データに相当する。また、説明を簡単にするため、「購入する」というクラスを正（+）のクラスと表し、「購入しない」というクラスを負（-）のクラスと表す場合がある。決定木では、クラス毎（例えば、「購入する（正）」、「購入しない（負）」というクラス毎）に個別データ数をまとめた情報をノードとする。例えば、図 2 に示すルートのノードでは、「する：9、しない：5」という情報をノードとしている。

#### 【0038】

分類器生成部 1 2 は、ステップ A 4 で生成された部分集合が与えられると、どの属性で最初にルートのノードを分割させるかを決定する。また、個別データに含まれている属性の個数を R 個とした場合、各属性を属性 1 ~ R とする。このとき、分類器生成部 1 2 は、属性 1 ~ R の各属性について、分割時の評価値を計算し、その評価値が最大の属性を、分割に最も適した属性として選択する。ここでは、分割前のノードのエントロピーと、分割後のエントロピーの差を評価値とする場合を例にするが、他の計算方法で評価値を求めてもよい。ノードのエントロピーは、クラスが正（+）の個別データの割合を q とし、クラスが負（-）の個別データの割合を 1 - q とすると、 $-q \log q - (1 - q) \log (1 - q)$  で表される。分割後のノードのエントロピーは、分割後の各ノードのエントロピーの加重平均である。

40

#### 【0039】

例えば、ステップ A 4 で生成された部分集合において、「正」が 9 データあり、「負」

50

が5データあるとするとルートのノードは「正：9，負：5」となる。この場合、正（+）のデータが9データあり、負（-）のデータが5データあるので、ルートのノードのエントロピーは、 $-(9/14) \times \log(9/14) - (5/14) \times \log(5/14) = 0.940$ となる。ただし、本例では $\log$ の底を2とする。

【0040】

分類器生成部12は、一つの属性でルートノードを分割して得られるノードを求める。すなわち、その属性の属性値毎に、正および負の個別データ数を表す情報（ノード）を生成する。例えば、その属性1のとり得る値が「0」または「1」であり、属性1の値が「0」のときには、正が5データあり、負が2データあるとし、属性1の値が「1」のときには、正が0データあり、負が7データあるとする。この場合、属性1の値が「0」か「1」かで分岐するノードとして、「正：5，負：2」というノードと、「正：0，負：7」というノードとを生成する。分類器生成部12は、分割後の各ノードのエントロピーを計算し、分割後の各ノードにおける正または負としてカウントされる個別データ数に応じて各ノードのエントロピーの加重平均を求める。上記の例では「正：5，負：2」というノードにおいても、「正：0，負：7」というノードにおいても個別データの総数は7であるので、加重平均を行う際の重み付け係数は各ノードでいずれも $(7/14)$ となる。従って、本例の場合、分類器生成部12は、分割後のエントロピーを以下のように計算する。

【0041】

$$(7/14) \times \{ -(5/7) \times \log(5/7) - (2/7) \times \log(2/7) \} + (7/14) \times \{ -(0/7) \times \log(0/7) - (7/7) \times \log(7/7) \} = 0.432$$

【0042】

ただし、 $\log$ の係数が0となる場合、その項の値は0とする。上記の例では、 $-(0/7) \times \log(0/7)$ の値は0としている。

【0043】

従って、本例の場合、分類器生成部12は、属性1で分割した場合の評価値を、 $0.940 - 0.432 = 0.508$ と計算する。

【0044】

分類器生成部12は、属性1だけでなく、他の属性についても同様に、その属性で分割したときの評価値を計算し、評価値が最大となる属性で分割すると決定する。このようにして、ルートノードを分割する属性を決定する。

【0045】

なお、上記の属性1の例では、属性1のとり得る値が「0」または「1」の二つだけである場合を示した。属性値が年齢であり、その値が20，21，22のように連続する値の場合には、どの属性値で分割させるのかも決める。この場合、分類器生成部12は、各属性値間の中間値をしきい値とし、各しきい値毎に、その「しきい値以下」および「そのしきい値より大」とに分割させた場合の評価値を求める。そして、評価値が最大となる場合を選択することによって、どの属性値で分割させるのかも決定する。例えば、属性値が20，21，22，・・・と連続する場合では、「20.5以下」および「20.5より大」で分割した場合の評価値、「21.5以下」および「21.5より大」で分割した場合の評価値等をそれぞれ計算し、評価値が最も高くなるように分割すればよい。

【0046】

分類器生成部12は、分割後の各ノードについても、上記と同様の処理を行い、次にどの属性で分割するのかを決定する処理を順次、繰り返す。また、分類器生成部12は、所定の条件が満たされたときには、ノードの分割を停止する。所定の条件とは、例えば、「ノードにおける個別データのクラスが全て同じになる」という条件や、「ノードにおける正または負としてカウントされる個別データ数が所定数（例えば2）以下になる」という条件を用いてよい。前者の条件を採用すると、ノードにおける個別データが全て正または負になると、そのノードの分割を継続しない。このように、分類器生成部12は、ルート

10

20

30

40

50

のノードから順次、分割を繰り返し、木構造の決定木を生成する。

【 0 0 4 7 】

また、分類器生成部 1 2 は、上記のように、木構造の決定木を生成した後、その決定木に対する枝刈りを行う。決定木において、分割されて生成された最終的なノードを葉と呼ぶ。ある葉に分類されたデータ数が  $D$  であるとする（すなわち、正または負としてカウントされる個別データ数が  $D$  であるとする）。この葉に分類された  $D$  データ中、 $E$  データが誤りであるとする。この仮定では、 $D$  回の試行中、誤りという事象を  $E$  回観測したとみなし、大きさ  $D$  の標本で、誤りという事象が起きる確率が  $r$  である二項分布と考えることができる。予め与えられた信頼度  $CF$  に対して、 $r$  の上限を  $U\_CF(E, D)$  と表すことにすると、 $D$  データでの誤りの発生する期待値は、 $D \times U\_CF(E, D)$  となる。分類器生成部 1 2 は、子のノードが全て葉である親のノードに対し、親における誤りの期待値（誤りの発生する期待値）と、子である葉の誤りの期待値の合計とを比較する。そして、子での期待値の合計の方が親の誤りの期待値よりも大きければ、分類器生成部 1 2 は、葉を縮退して、その親を葉とする。分類器生成部 1 2 は、この処理を順次繰り返すことで、決定木全体の葉の枝刈りを行う。

10

【 0 0 4 8 】

葉を縮退する場合、分類器生成部 1 2 は、葉を削除して、その削除した葉の親のノードを葉とすればよい。例えば、図 4 に例示する決定木において、「湿度」という属性の値に応じて分割したノードを縮退する場合、分類器生成部 1 2 は、湿度の属性値が 70% 以下となっている個別データ数を表すノード「する：2，しない：0」と、湿度の属性値が 70% より高くなっている個別データ数を表すノード「する：0，しない：3」とを削除して、その 2 つのノードの親ノード「する：2，しない：3」を葉とすればよい。

20

【 0 0 4 9 】

分類器として決定木を生成する場合、例えば、上記のように、決定木を定めて枝刈りを行うことで、決定木を生成すればよい。

【 0 0 5 0 】

分類器生成部 1 2 は、ステップ A 4 で生成された部分集合に属するデータのうち、特定のクラスのデータを重く加重してから決定木を生成してもよい。ここで、加重とは、データに対する重みを設定することである。例えば、部分集合に属するデータのうち、予め定められた特定のクラスの個別データの数が数倍になるように、そのクラスの個別データの複製を作成することで、加重を行ってもよい。なお、特定のクラスの個別データ数を何倍にするかは予め定めておけばよい。そして、特定のクラスの個別データ数を増やすように、そのクラスの個別データを複製した後に、分類器を生成してもよい。加重の対象とするクラスおよび加重量（例えば、データを何倍に増やすか等）は、キーボード等の入力装置（図示せず）を介してデータ類似度計算システムのユーザにより入力されてもよいし、他の様態で指定されてもよい。

30

【 0 0 5 1 】

ステップ A 5 で分類器が生成されると、自己評価部 1 3 は、その分類器を用いて、ステップ A 3 で生成された部分集合に属する個々のデータのクラスを判定する（ステップ A 6）。ここでも、分類器が決定木である場合を例にして説明する。分類器となる決定木を生成した場合、その決定木と、既知の属性（部分集合に属する個別データの属性）とから、クラスを予測する。このとき、自己評価部 1 3 は、決定木のルートのノードを起点として、ノードを分割する際に用いた属性に関して個別データの属性値を参照し、その属性値に応じて子ノードを辿る。自己評価部 1 3 は、子ノードを辿っていき、葉のノードまで辿ったならば、葉のノードでカウント数の多い方のクラスを、分類器を用いた判定の結果とすればよい。

40

【 0 0 5 2 】

図 5 は、分類器から判定されたクラスの例を示す説明図である。図 5 に示す各属性の属性値は、図 2 に示す各属性値と同一である。また、図 5 では、図 4 に例示する決定木（分類器）を用いて、その属性値からクラスを判定した結果を示している。例えば、個別デー

50

タにおける属性「天気」の属性値が「雨」であるとする。すると、ルートのノードから、「する：3、しない：2」という子ノードを辿る（図4参照）。そのノードが葉であるので、自己評価部13は、「する：3、しない：2」というカウント数により、「購入しない」というクラスであると判定する。このように判定した結果、図5に示す6番目および最後の個別データのクラスの判定結果は、図2に示す元のクラスと異なっている（図2、図5参照）。

【0053】

ステップA6で部分集合に属する各個別データのクラスが判定された後、類似度算出部14は、データの集合に属する個別データを順番に指定するための第1の変数*i*を1に初期化する（ステップA7）。続いて、類似度算出部14は、変数*i*の値が個別データの総数*N*以下であるか否かを判定する（ステップA8）。変数*i*の値が*N*を超えていれば（ステップA8におけるNo）、変数*t*を1インクリメントし（ステップA8）、その後、データ類似度計算システムはステップA3以降の処理を再度繰り返す。

10

【0054】

変数*i*の値が*N*以下であるならば（ステップA8におけるYes）、類似度算出部14は、変数*i*によって定まる*i*番目の個別データが、直近のステップA4で生成された部分集合に含まれているか否かを判定する（ステップA10）。*i*番目のデータが部分集合に含まれていないならば（ステップA10におけるNo）、類似度算出部14は、変数*i*を1インクリメントし（ステップA11）、ステップA8以降の処理を再度繰り返す。

20

【0055】

*i*番目のデータが部分集合に含まれているならば（ステップA10におけるYes）、データの集合に属する個別データを順番に指定するための第2の変数*j*の値を*i*+1に設定する（ステップA12）。すなわち、*j*に*i*+1を代入する。

【0056】

ステップA12の後、類似度算出部14は、変数*j*の値が*N*以下であるか否かを判定する（ステップA13）。変数*j*の値が*N*を超えていれば（ステップA13におけるNo）、類似度算出部14は、変数*i*を1インクリメントし（ステップA11）、ステップA8以降の処理を再度繰り返す。

【0057】

変数*j*の値が*N*以下であるならば（ステップA13におけるYes）、類似度算出部14は、変数*j*によって定まる*j*番目の個別データが、直近のステップA4で生成された部分集合に含まれているか否かを判定する（ステップA14）。*j*番目の個別データが部分集合に含まれていないならば（ステップA14におけるNo）、類似度算出部14は、変数*j*を1インクリメントし（ステップA17）、ステップA13以降の処理を繰り返す。

30

【0058】

*j*番目の個別データが部分集合に含まれているならば（ステップA14におけるYes）、類似度算出部14は、*i*番目の個別データの属性から判定されたクラス（PredC<sub>*i*</sub>と記す）と、*j*番目の個別データの属性から判定されたクラス（PredC<sub>*j*</sub>と記す）とを比較し、両者が同じクラスであるか否かを判定する（ステップA15）。PredC<sub>*i*</sub>およびPredC<sub>*j*</sub>は、直近のステップA6で、*i*番目および*j*番目の個別データについて、分類器を用いて属性から判定されたクラスである。PredC<sub>*i*</sub>とPredC<sub>*j*</sub>とが同じクラスであるということは、*i*番目の個別データおよび*j*番目の個別データが、同一のクラスと判定されたデータ同士であるということの意味する。

40

【0059】

PredC<sub>*i*</sub>とPredC<sub>*j*</sub>とが異なっていれば（ステップA15におけるNo）、類似度算出部14は、変数*j*を1インクリメントし（ステップA17）、ステップA13以降の処理を繰り返す。

【0060】

PredC<sub>*i*</sub>とPredC<sub>*j*</sub>とが同じクラスであれば（ステップA15におけるYes）、*i*番目の個別データと*j*番目の個別データの類似度Sim(*i*, *j*)に所定数を加算

50

する(ステップA16)。本例では、この所定数を1とし、ステップA16で $Sim(i, j)$ に1を加算する場合を例にする。ステップA16の後、類似度算出部14は、変数 $j$ を1インクリメントし(ステップA17)、ステップA13以降の処理を繰り返す。

【0061】

上記のように、1つの部分集合が生成されたとき、その部分集合に属している1対の個別データであって、 $PredC_i$ と $PredC_j$ とが同じクラスとなっている個別データ間の類似度 $Sim(i, j)$ の値を1増加させていく。また、一方あるいは両方が部分集合に属していない個別データ同士の場合には、ステップA16に移行しないので、その個別データ間の類似度は増加しない。また、部分集合に属している1対の個別データであっても、属性から判定されたクラスが異なる場合にも、ステップA16に移行せず、その個別データ間の類似度は増加しない。

10

【0062】

そして、データ類似度計算システムは、変数 $t$ をインクリメントして(ステップA8)、部分集合を生成すると(ステップA4)、ステップA5以降の処理を実行する。このとき、各個別データの組について一律にステップA16の処理を実行するわけではない。個別データの重複を許して部分集合の生成を複数回生成したときに、ステップA16が行われた回数が多い個別データの組については、 $Sim(i, j)$ の値も大きくなる。一方、ステップA16が行われた回数の少ない個別データの組については、 $Sim(i, j)$ の初期値からの増加量は少ない。部分集合の生成回数 $t$ が上限値 $T$ に達したときにおける各データの対毎の $Sim(i, j)$ が、対をなす個別データ間の類似度となる。

20

【0063】

また、上記の第1の実施の形態では、二つのデータがいずれも部分集合に属し、その二つのデータに対して判定されたクラスが同一であるという条件が満たされたときに、その二つのデータ類似度 $Sim(i, j)$ に対して所定値を加算する。この条件だけでなく、さらに、二つのデータに対して判定されたクラスが特定のクラス(例えば、「購入する」)であるという条件を満たしている場合に、二つのデータ類似度 $Sim(i, j)$ に対して所定値を加算し、他の場合には、 $Sim(i, j)$ への加算を行わなくてもよい。この場合、二つのデータに対するクラスの判定結果が同一であっても、そのクラスが特定のクラスでない場合、類似度に対する加算を行わない。

【0064】

30

上記のように類似度を求める場合、類似度算出部14は、ステップA15において、 $PredC_i$ と $PredC_j$ とが同一であり、かつ、 $PredC_i$ および $PredC_j$ が特定のクラス(例えば、「購入する」)であるか否かを判定すればよい。そして、この条件を満たしている場合に、ステップA16を行い、満たしていない場合には、ステップA17に移行すればよい。

【0065】

次に、第1の実施形態の効果について説明する。

第1の実施形態によれば、与えられたデータの集合に属するデータ同士の組に対して、それぞれ類似度の初期値を定める。そして、そのデータの集合から部分集合を生成し、その部分集合に含まれる各データに定められている属性(例えば、顧客の特徴や販売条件)およびクラスから、分類器を生成する。さらに、分類器を用いて、その属性からクラスを判定し、部分集合に含まれているデータ同士であって、判定されたクラスが同じデータの類似度に所定値を加算する。上記の部分集合生成以後の処理を複数回繰り返すことで、各データ同士の類似度を決定する。本願発明では、このような処理によって、自動的に、各データ間の類似度を算出することができる。

40

【0066】

また、本実施形態において、部分集合生成時(ステップA4)では、個別データを予め定められた数 $M$ だけ選択すればよく、既に生成された部分集合に属している個別データを重複を許して選択する。例えば、ランダムサンプリングを行ってもよい。このように、 $M$ 個の個別データを選択することを複数回行って、類似度を生成する。また、類似度が求ま

50

れば、類似度に基づいて個別データを分類することができる。特許文献1に記載されたシステムのように、分析者が手動で顧客の分類を行う場合には、顧客の分類に対する特別なスキルが必要となったり、試行錯誤して分析を行って手間がかかったりすることがあるが、上記のように類似度を求めれば、システムのユーザに特別なスキルがなくても、ユーザの手間をかけずに、分類に用いるための類似度を算出することができる。

【0067】

また、各データが個々の顧客に対応し、データの属性が顧客の特徴や販売条件であり、データに「購入する」または「購入しない」というクラスが与えられている場合、商品またはサービスに対する顧客の行動（購入したか否か）に基づいて、データの類似度を求めることができる。また、データが顧客に対応しているため、顧客同士の類似度ということもできる。

10

【0068】

以下に示す効果の説明では、データの属性が顧客の特徴であり、データと顧客が対応していて、データのクラスが顧客の購入行動（例えば「購入する」または「購入しない」）である場合を例にして説明する。

【0069】

重複を許したランダムサンプリングにより部分集合を生成した場合、個々のデータ（本例では個々の顧客）は、部分集合に $M/N$ の確率で選択されることとなる。すなわち、 $M/N$ の確率で部分集合に含まれることになる。他のデータ（他の顧客）も、同様の確率で部分集合に含まれる。データを組み合わせて得られる各データの組は、同じ確率で部分集合に属することとなる。従って、部分集合を作成する回数 $T$ が十分大きく、仮に、元のクラスとステップA6で判定されたクラスが全て同一であれば、データ間の類似度はほぼ同じとなり、クラスが同一のデータを細かく分類することはできない。例えば、「購入する」というクラスのデータを細かく分類したり、「購入しない」というクラスのデータを細かく分類したりすることはできない。また、「購入する」というクラスのデータ（商品等を購入した顧客）と、「購入しない」というクラスのデータ（商品等を購入しなかった顧客）との類似度は0となる。これに対し、本発明の第1の実施形態では、部分集合に属する各個別データから分類器を生成し、その分類器を用いて、個別データの属性からクラスを判定することで、元のクラスと、判定後のクラスとがことなる場合を生じさせる機会を積極的に設けている。そして、判定後のクラスを用いて、個別データの類似度を求めることで、 $i$ 番目のデータと $j$ 番目のデータからなる各組に対して、それぞれ類似度 $S_{im}(i, j)$ を定めることができる。

20

30

【0070】

また、本発明では、互いに関連のある属性が多く存在しても、それらの関連ある属性の影響を強く受けることなく類似度を算出できる。例えば、図12に示すデータ番号1~10の10個のデータが与えられているとする。図12に示すデータ番号{1, 2, 7, 9, 10}のデータを部分集合とすると、例えば、分類器として「属性aの値が1ならばクラスは『購入する』であり、属性aの値が2ならばクラスは『購入しない』である」という決定木が生成され、データ番号{1, 2}のデータ間の類似度や、データ番号{7, 9, 10}の各データ間の類似度が加算される。また、例えば、データ番号{1, 2, 6, 7, 9, 10}のデータを部分集合とする。この場合、例えば、分類器として「属性bの値が1ならばクラスは『購入する』であり、属性bの値が2ならばクラスは『購入しない』である」という決定木が生成され、データ番号{1, 2, 6}の各データ間の類似度や、データ番号{7, 9, 10}の各データ間の類似度が加算される。

40

【0071】

一般に、決定木には複数の属性が現れるが、属性に互いに関連性がある場合、決定木に現れる属性の数は減る。上記のデータ番号{1, 2, 6, 7, 9, 10}を部分集合とする場合を例にすると、属性bを用いて個別データを2つに分類した場合と、属性c、d、あるいはeを用いて個別データを2つに分類した場合とで、個別データの分類のされ方が全く同一になる。よって、属性bで分割した後、属性c、d、あるいはeで分割すること

50

はない。具体例を挙げると、データ番号 { 1 , 2 , 6 , 7 , 9 , 10 } のデータを属性 b で分割すると、属性 b の値が 1 のノードにはデータ番号 { 1 , 2 , 6 } のデータが含まれ、属性 b の値が 2 のノードにはデータ番号 { 7 , 9 , 10 } が含まれる。データ番号 { 1 , 2 , 6 } のノードを c で分割する場合、そのデータ番号 { 1 , 2 , 6 } のデータは全て属性 c の値が 1 となり、属性 c の値が 2 となるデータはない。属性 b の値が 2 のノードでも、属性 c の値が一方に偏っている。属性 b での分割後、属性 d や e で分割する場合も同様である。

**【 0 0 7 2 】**

このように、属性に互いに関連性がある場合、決定木に現れる属性の数が減り、関連のある属性の影響を強く受けることがなくなる。上記のデータ番号 { 1 , 2 , 6 , 7 , 9 , 10 } を部分集合とする例では、例えば、決定木には属性 b が現れるが、属性 b と同時に他の属性 c、d、e は決定木に現れず、属性 b に関連する属性 c、d、e に影響されることがない。そのため、データに関連ある属性が多く存在していても、それらの関連のある属性の影響を強く受けることなく、類似度を求めることができる。

10

**【 0 0 7 3 】**

また、二つのデータがいずれも部分集合に属し、その二つのデータに対して判定されたクラスが同一であるという条件だけでなく、さらに、その二つのデータに対して判定されたクラスが特定のクラス（例えば、「購入する」）であるという条件を満たしているときにのみ、その二つのデータの類似度  $Sim(i, j)$  への加算（ステップ A 16）を行う場合、以下の効果が得られる。すなわち、特定のクラスのデータの持つ特徴の違いを重視して類似度を算出できる。ここでも、データと顧客とが対応している場合を例にして説明する。ステップ A 16 に移行する条件として、二つのデータに対して判定されたクラスが特定のクラスであるという条件も加えると、特定のクラスの顧客が持つ特徴の違いを重視して類似度を算出できる。例えば、商品またはサービスを購入した顧客が複数の特徴のいずれかを持っていて、まだその商品またはサービスを購入していないが、既に購入済みの顧客と特徴が近い顧客（購入の見込みのある顧客）が複数の特徴のうちのいずれかを持っていそうな場合がある。そのような場合において、購入していない顧客との類似度は重視せず、購入した顧客との類似度を重視したいと分析者が考える場合がある。そのような場合、ステップ A 6 で「購入する」という特定のクラスであると判定された顧客のデータと、ステップ A 6 で「購入する」と判定された「購入の見込みのある顧客」のデータとの類似度について加算して、「購入する」というクラスの顧客が持つ特徴の違いを重視して類似度を算出することができる。

20

30

**【 0 0 7 4 】**

また、分類器生成部 12 は、分類器を作成するときに、ステップ A 4 で生成された部分集合に属するデータのうち、予め与えられたクラスが特定のクラスとなっているデータに対して加重を行ってから、分類器を作成してもよい。例えば、部分集合に属するデータのうち、予め与えられたクラスが「購入する」となっているクラスのデータのデータ数が、指定された分だけ増えるように、そのデータの複製を作成して、その複製したデータも用いて分類器を生成してもよい。そのような分類器を用いて類似度を算出すれば、その類似度に基づいてデータを分類する場合、特定のクラス以外のクラスのデータが持つ特徴の違いを重視してクラスタ（グループ）に分類することができる。例えば、「購入する」というクラスの顧客（データ）と、「購入しない」というクラスの顧客の数が同じくらいであるとする。データに加重を行わない場合には、いずれのクラスのデータも混ざったクラスタが生成されがちになる。特定のクラスのデータを重く加重すると、特定のクラスのデータの間の類似度がどの組み合わせでも高くなることで 1 つのクラスタとなる。一方、特定のクラス以外のクラスのデータは、そのように類似度が高くなることはなく、1 つのクラスタにまとめられずに、複数のクラスタに分類することができる。なお、類似度を用いてデータをクラスタ（グループ）に分類する処理については、第 2 の実施形態で説明する。

40

**【 0 0 7 5 】**

実施形態 2 .

50

図6は、本発明の第2の実施形態の例を示すブロック図である。第2の実施形態のデータ類似度計算システム20は、部分集合生成部11と、分類器生成部12と、自己評価部13と、類似度算出部14と、類似度クラスタリング部21とを備える。第1の実施形態と同様の構成要素については、図1と同一の符号を付し、詳細な説明を省略する。第2の実施形態のデータ類似度計算システム20は、第1の実施形態と同様に、データの集合に属する個別データ間の類似度を計算した後、その類似度を用いて個別データを分類する。従って、第2の実施形態のデータ類似度計算システム20は、分類システムと称することができる。以下、本実施形態において、個別データのグループをクラスタと記す。

【0076】

類似度クラスタリング部21は、類似度算出部14が求めた個別データ間の類似度に基づいて、データの集合に属する個別データを複数のクラスタに分類する。類似度クラスタリング部21には、目標とするクラスタ数が入力され、類似度クラスタリング部21は、そのクラスタ数になるように個別データを分類する。目標とするクラスタ数は、データ類似度計算システム20に設けられるキーボード等の入力装置(図示せず)を介して、データ類似度計算システムのユーザによって入力されてもよい。あるいは、他の態様でクラスタ数が指定されてもよい。

10

【0077】

部分集合生成部11、分類器生成部12、自己評価部13、類似度算出部14および類似度クラスタリング部21は、例えば、プログラム(データ類似度計算プログラム)に従って動作するCPUによって実現される。その場合、CPUがプログラムに従って、部分集合生成部11、分類器生成部12、自己評価部13、類似度算出部14および類似度クラスタリング部21として動作すればよい。

20

【0078】

次に、第2の実施形態の動作について説明する。

データ集合に属する各個別データ間の類似度を求める処理は、第1の実施形態と同様である。部分集合生成部11、分類器生成部12、自己評価部13および類似度算出部14がそれぞれ第1の実施形態と同様に動作し、例えば、図3に示す処理を行って、 $i, j$ の組毎に、個別データ間の類似度 $Sim(i, j)$ を求めればよい。

【0079】

それぞれの個別データ間の類似度 $Sim(i, j)$ が算出された後、類似度クラスタリング部21は、個別データの各組における類似度を用いて、個別データをクラスタリングする(すなわち分類する)。

30

【0080】

類似度が与えられたときのクラスタリング方法には、様々な方法がある。例えば、階層的クラスタリング法として最短距離法、最長距離法、郡平均法、ワード法等があり、非階層的クラスタリング法としてK平均法等がある。類似度クラスタリング部21は、いずれの方法で個別データを分類してもよい。また、どの方法でクラスタリングを行うかの指定が、キーボード等の入力装置を介してデータ類似度計算システムのユーザにより入力されてもよい。あるいは、他の態様で指定されてもよい。

【0081】

40

以下、最短距離法によって個別データをクラスタリングする場合を例にして、個別データのクラスタリング処理を説明する。図7は、最短距離法によって個別データをクラスタリングする処理経過の例を示すフローチャートである。なお、以下の説明では、目標とするクラスタ数をK個とする。目標とするクラスタ数Kは、例えば、キーボード等の入力装置を介して、予め類似度クラスタリング部21に入力される。

【0082】

類似度クラスタリング部21は、データの集合に属する個別データをそれぞれ1個だけ含むクラスタを、各個別データ毎に定める(ステップB1)。従って、類似度クラスタリング部21は、ステップB1において、データの集合に属する個別データの総数Nと等しいN個のクラスタを定めることになる。このN個のクラスタに属する個別データは、クラ

50



スタ毎に異なっている。ステップ B 1 の後、ステップ B 2 に移行する。

【 0 0 8 3 】

ステップ B 2 では、類似度クラスタリング部 2 1 は、二つのクラスタからなるクラスタの各組についてそれぞれ、クラスタ間の類似度を求め、最もクラスタ間の類似度が高い二つのクラスタを特定し、その二つのクラスタを一つのクラスタに併合する（ステップ B 2）。ステップ B 2 に移行した時点でのクラスタ数を L とすると、L 個のクラスタから二つのクラスタを取り出す  $L C_2$  個の組毎に、クラスタ間の類似度を求め、最もクラスタ間の類似度が高い二つのクラスタを一つのクラスタに併合すればよい。なお、最初にステップ B 1 からステップ B 2 に移行したときには、 $L = N$  である。

【 0 0 8 4 】

また、二つのクラスタからなる組において、その二つのクラスタ間の類似度を定める方法の例を以下に示す。類似度を求める対象となる二つのクラスタの一方を  $C_1$  と記し、他方を  $C_2$  と記す。また、その二つのクラスタ  $C_1, C_2$  間の類似度を  $S(C_1, C_2)$  と記す。類似度クラスタリング部 2 1 は、例えば、 $C_1$  に属する個別データと、 $C_2$  に属する個別データとの各組み合わせにおける個別データ間の類似度のうち、最大値を  $C_1, C_2$  間の類似度を  $S(C_1, C_2)$  と定めればよい。すなわち、 $C_1$  から取り出した一つの個別データを  $x_1$  とし、 $C_2$  から取り出した一つの個別データを  $x_2$  とし、 $x_1, x_2$  の類似度を  $S(x_1, x_2)$  とすると、類似度クラスタリング部 2 1 は、式 (1) に示すように  $S(x_1, x_2)$  の最大値を  $S(C_1, C_2)$  と定めればよい。

【 0 0 8 5 】

【 数 1 】

$$S(C_1, C_2) = \max_{x_1 \in C_1, x_2 \in C_2} S(x_1, x_2) \quad \text{式 (1)}$$

【 0 0 8 6 】

また、二つのクラスタを一つのクラスタに併合するとは、二つのクラスタに属する各個別データを一つのクラスタにまとめることである。二つのクラスタを一つのクラスタに併合することにより、クラスタの総数が 1 つ減少する。

【 0 0 8 7 】

ステップ B 2 で二つのクラスタを一つのクラスタに併合した後、類似度クラスタリング部 2 1 は、クラスタ数が目標数 K になったか否かを判定する（ステップ B 3）。クラスタ数が目標数 K まで減っていなければ（ステップ B 3 における No）、ステップ B 2 以降の処理を繰り返す。クラスタ数が目標数 K となっていれば（ステップ B 3 における Yes）、K 個のクラスタが得られ、個別データが K 個に分類されているので、処理を終了する。

【 0 0 8 8 】

第 2 の実施形態によれば、データの集合に含まれる個別データを目標数のクラスタに分類することができる。

【 0 0 8 9 】

例えば、データの属性が顧客の特徴であり、データと顧客が対応していて、データのクラスが顧客の購入行動（「購入する」または「購入しない」等）を表している場合、商品またはサービスに対する顧客の購入行動に応じて、自動的に顧客を分類することができる。

【 0 0 9 0 】

また、図 7 に例示するように、最初に各個別データをそれぞれ別々のクラスタに振り分け、クラスタの数が目標数となるまでクラスタを併合させていけば、目標数のクラスタに個別データを分類することができる。すなわち、特定のクラスタに個別データが集まってしまい個別データのクラスタ数が目標数に達しないということを防止することができる。

【 0 0 9 1 】

実施形態 3 .

図8は、本発明の第3の実施形態の例を示すブロック図である。第3の実施形態のデータ類似度計算システム30は、部分集合生成部11と、分類器生成部12と、自己評価部13と、類似度算出部14と、類似度クラスタリング部21と、属性データ分類部31と、関連性算出部32とを備える。第1の実施形態や第2の実施形態と同様の構成要素については、図1、図6と同一の符号を付し、詳細な説明を省略する。第3の実施形態のデータ類似度計算システム30は、第2の実施形態と同様に類似度を用いて個別データを分類するので、分類システムと称することができる。なお、本実施の形態においても、類似度クラスタリング部21によって分類された個別データのグループをクラスタと記す。また、第3の実施形態のデータ類似度計算システム30は、類似度算出に対する属性の関連度を算出する。類似度算出に対する属性の関連度とは、類似度算出に対する属性の関連性の度合い（換言すれば、類似度算出に対して属性が影響を与える度合い）を示す数値である。

10

【0092】

属性データ分類部31には、属性およびその属性値に基づく分類方法を指定する情報（以下、分類方法指定情報と記す）が入力され、属性データ分類部31は、指定された分類方法に従って、個別データを分類する。分類方法指定情報によって、属性データ分類部31による個別データの分類数も定まる。分類方法指定情報の例として、「属性名Aの属性値がThres以上である個別データをグループ1に分類し、属性名Aの属性値がThres未満である個別データをグループ2に分類する。」などの情報が挙げられる。この場合、個別データは、二つのグループ1, 2に分類されることになる。分類方法指定情報は、例えば、キーボード等の入力装置（図示せず）を介して属性データ分類部31に入力されるが、他の態様で入力されてもよい。

20

【0093】

関連性算出部32は、類似度クラスタリング部21によって分類されたクラスタ（グループ）と、属性データ分類部31によって分類されたグループとの関係に基づいて、類似度算出に対する属性の関連度（以下、単に属性の関連度と記す）を求める。この属性は、分類方法指定情報で指定された属性である。

【0094】

部分集合生成部11、分類器生成部12、自己評価部13、類似度算出部14、類似度クラスタリング部21、属性データ分類部31および関連性算出部32は、例えば、プログラム（データ類似度計算プログラム）に従って動作するCPUによって実現される。その場合、CPUがプログラムに従って、部分集合生成部11、分類器生成部12、自己評価部13、類似度算出部14、類似度クラスタリング部21、属性データ分類部31および関連性算出部32として動作する。

30

【0095】

次に、第3の実施形態の動作について説明する。

類似度クラスタリング部21が個別データを分類するまでの動作は、第2の実施形態と同様である。部分集合生成部11、分類器生成部12、自己評価部13および類似度算出部14は、第1および第2の実施形態と同様に動作し、例えば、図3に示す処理を行って、i, jの組毎に、個別データ間の類似度Sim(i, j)を求めればよい。その後、類似度クラスタリング部21は、第2の実施形態と同様に、指定されたクラスタ数に個別データを分類する。

40

【0096】

図9は、類似度クラスタリング部21による分類後に属性の関連度を求める動作の例を示すフローチャートである。ここでは、類似度クラスタリング部21によってグループ分けされるクラスタ数と、属性データ分類部31によってグループ分けされるグループ数が等しい場合を例に説明する。本例では、類似度クラスタリング部21がデータの集合に属する個別データを二つのクラスタに分類するものとする。また、属性データ分類部31には、例えば、「属性名Aの属性値がThres以上である個別データをグループ1に分類し、属性名Aの属性値がThres未満である個別データをグループ2に分類する。」と

50

いう分類方法指定情報が入力され、属性データ分類部 3 1 が個別データを 2 つのグループに分類するものとする。

【 0 0 9 7 】

また、以下の説明において、変数  $k$  は、類似度クラスタリング部 2 1 によってグループ分けされるクラスタを指定するための変数であり、本例では、 $k$  の取り得る値は 1 または 2 である。また、変数  $l$  は、属性データ分類部 3 1 によってグループ分けされたグループを指定するための変数であり、本例では、 $l$  の取り得る値は 1 または 2 である。

【 0 0 9 8 】

属性データ分類部 3 1 は、まず、 $k = 1, 2$ 、 $l = 1, 2$  によって決まる  $k$  と  $l$  の組み合わせを順次定め、 $k, l$  の組み合わせと一対一に対応する変数  $N[k][l]$  を全て 0 に初期化する (ステップ C 1)。  $N[k][l]$  は、 $k$  番目のクラスタに属している個別データであって、 $l$  番目のグループにも属している個別データの数を表す。

10

【 0 0 9 9 】

属性データ分類部 3 1 は、個別データを指定するための変数  $i$  を 1 に初期化する (ステップ C 2)。

【 0 1 0 0 】

次に、属性データ分類部 3 1 は、変数  $k$  の値を、個別データ  $i$  (すなわち、 $i$  番目の個別データ) が属するクラスタのクラスタ番号とする (ステップ C 3)。なお、例えば類似度クラスタリング部 2 1 は、分類の結果得られた  $K$  個のクラスタに対して、クラスタを識別するための番号  $1 \sim K$  を割り当てる。個々のクラスタに割り当てられた番号がクラスタ番号である。ステップ C 3 では、変数  $k$  に、データ  $i$  が属するクラスタのクラスタ番号を代入すればよい。本例ではクラスタ数は 2 であるので、クラスタ番号は 1 または 2 である。

20

【 0 1 0 1 】

次に、属性データ分類部 3 1 は、個別データ  $i$  の属性値 (分類方法指定情報で指定された属性の属性値) に応じて、個別データ  $i$  をどのグループに含めるかを判定する (ステップ C 4)。すなわち、分類方法指定情報に従って、個別データ  $i$  に対する分類を行う。本例では、「属性名  $A$  の属性値が  $Thres$  以上である個別データをグループ 1 に分類し、属性名  $A$  の属性値が  $Thres$  未満である個別データをグループ 2 に分類する。」という分類方法指定情報に従って、個別データ  $i$  の属性  $A$  の属性値が  $Thres$  以上であるか否かを判定する。個別データ  $i$  の属性  $A$  の属性値が  $Thres$  以上であるならば (すなわち、個別データ  $i$  をグループ 1 に分類すると判定したならば)、変数  $l$  の値を 1 に設定する (ステップ C 5)。また、個別データ  $i$  の属性  $A$  の属性値が  $Thres$  未満であるならば (すなわち、個別データ  $i$  をグループ 2 に分類すると判定したならば)、変数  $l$  の値を 2 に設定する (ステップ C 6)。

30

【 0 1 0 2 】

ステップ C 5 またはステップ C 6 の後、ステップ C 3 で定められた  $k$  と、ステップ C 5 またはステップ C 6 で定められた  $l$  との組み合わせに対応する  $N[k][l]$  の値を 1 インクリメントする (ステップ C 7)。例えば、ステップ C 5 からステップ C 7 に移り、 $N[k][l]$  を 1 インクリメントした場合、 $k$  番目のクラスタに属している個別データであって、 $l$  番目のグループに属している個別データを一つカウントして、そのカウント値を 1 増加させたことになる。また、ステップ C 6 からステップ C 7 に移り、 $N[k][2]$  を 1 インクリメントした場合、 $k$  番目のクラスタに属している個別データであって、2 番目のグループに属している個別データを一つカウントして、そのカウント値を 1 増加させたことになる。

40

【 0 1 0 3 】

次に、属性データ分類部 3 1 は、変数  $i$  の値を 1 インクリメントする (ステップ C 8)。そして、属性データ分類部 3 1 は、変数  $i$  の値が個別データの総数  $N$  以下であるか否かを判定し (ステップ C 9)、 $i$  の値が  $N$  以下であれば (ステップ C 9 における  $Yes$ )、ステップ C 3 以降の処理を再度行う。従って、1 番目から  $N$  番目までの各個別データに対

50

してステップC 3以降の処理を行うことになり、k番目のクラスタに属している個別データであって、1番目のグループにも属している個別データの数N[k][1]が、k, 1の組み合わせ毎に求められる。

【0104】

iの値がNを超えていれば(ステップC 9におけるNo)、ステップC 10に移行する。

【0105】

ステップC 10において、関連性算出部32は、N[k][1]を用いて、類似度クラスタリング部による分類と属性データ分類部による分類との独立性を検定し、p値を算出する(ステップC 10)。このp値を、分類方法指定情報で指定された属性の関連度とすることができる。

10

【0106】

関連性算出部32は、例えば、類似度クラスタリング部21が行った分類と属性データ分類部31が行った分類とが独立であるという仮説によりp値を求める。この場合、関連性算出部32は、まずk, 1の各組み合わせに関して、N[k][1]の期待値(E[k][1]と記す)を計算する。上記の仮説のもとでは、関連性算出部32は、以下に示す式(2)の計算によって各N[k][1]の期待値E[k][1]を求めればよい。

【0107】

【数2】

20

$$E[k][l] = \frac{\sum_{g=1}^2 N[k][g] \times \sum_{h=1}^2 N[h][l]}{N} \quad \text{式(2)}$$

【0108】

そして、N[k][1]およびE[k][1]を用いて、以下に示す式(3)の計算によりχ<sub>0</sub><sup>2</sup>を計算すると、χ<sub>0</sub><sup>2</sup>はχ<sup>2</sup>分布に従う。

【0109】

【数3】

30

$$\chi_0^2 = \sum_{k=1}^2 \sum_{l=1}^2 \frac{(N[k][l] - E[k][l])^2}{E[k][l]} \quad \text{式(3)}$$

【0110】

そこで、関連性算出部32は、式(3)の計算によりχ<sub>0</sub><sup>2</sup>を求め、χ<sup>2</sup>分布表を参照してp値を決定すればよい。このp値が、分類方法指定情報で指定された属性の関連度である。なお、χ<sup>2</sup>分布表は、例えば、予めデータ類似度計算システム30が備える記憶装置(図示せず)に記憶させておけばよい。

40

【0111】

χ<sup>2</sup>分布表からp値を決定する際に用いる自由度は、類似度クラスタリング部21によって分類されるクラスタ数-1である。従って、本例では、クラスタ数=2であるので、自由度は、2-1=1とすればよい。

【0112】

図10は、χ<sup>2</sup>分布表の例を示す説明図である。図10では、自由度をdfで表し、df=1, 2, 3の場合を例示しているが、df=4以上の場合も含めておく。図10に示すχ<sup>2</sup>分布表の最上段の値「0.7」、「0.5」、「0.3」、「0.2」、「0.1」等はp値である。例えば、式(3)でχ<sub>0</sub><sup>2</sup>を計算した結果、関連性算出部32は、χ<sub>0</sub><sup>2</sup>

50

1.07であったとする。本例では、 $\chi_0^2$  は、自由度1の $\chi^2$ 分布に従うので、自由度1における $\chi_0^2 = 1.07$ に応じたp値“0.3”を $\chi^2$ 分布表から特定し、そのp値“0.3”を、類似度算出に対する属性Aの関連度とすればよい。

【0113】

以上のように、属性Aの関連度を求めることができる。他の属性の関連度も、分類方法指定情報を入力して求めることができる。

【0114】

また、関連性算出部32は、例えば、 $\chi_0^2$ を計算した後に以下に示す式(4)の計算を行って、クラメールの関連係数を求め、そのクラメールの関連係数を、類似度算出に対する属性の関連度としてもよい。

【0115】

【数4】

$$\sqrt{\frac{\chi_0^2}{N}}$$

式(4)

【0116】

クラメールの連関係数は、0から1までの数値であり、1に近いほど強く関連していることを示す。

【0117】

また、関連性算出部32が類似度算出に対する属性の関連度を求めた後、その属性の関連度を、ディスプレイ装置またはプリンタ装置等の出力装置(図示せず)に出力させてから、終了してもよい。ユーザが入力装置(図示せず)を介して、属性を指定した分類方法指定情報を入力し、データ類似度計算システム30がその属性の関連度を求めてもよい。そして、関連性算出部32は、そのようにして求めた各属性の関連度をディスプレイ装置またはプリンタ装置等の出力装置(図示せず)に出力させる際、表形式やグラフ形式で出力させてもよい。

【0118】

次に、第3の実施形態の効果について説明する。

本発明によれば、類似度によるクラスタリングと属性値との関連の度合いを調べることができる。例えば、商品またはサービスに対する顧客の行動に応じて算出した類似度によるクラスタリングと属性値との関連を調べることができる。具体例を挙げると、例えば、類似度によるクラスタリングと男女(性別)との関連性が高いかどうかを調べることができる。商品またはサービスに対する顧客の行動に応じて、顧客を分類する際に関連のある顧客の特徴や販売条件を抽出することができる。これにより、商品またはサービスの今後ターゲットとすべき顧客等を分析することができる。

【0119】

次に、本発明の概要について説明する。図11は、本発明の概要を示すブロック図である。本発明のデータ類似度計算システム80は、部分集合生成手段81と、分類器生成手段82と、クラス判定手段83と、類似度算出手段84とを備える。

【0120】

部分集合生成手段81(例えば、部分集合生成部11)は、データの特徴を示す属性とデータの類別を示すクラスとを含むデータの集合から、データの重複を許してその集合の部分集合を複数回生成する。

【0121】

分類器生成手段82(例えば、分類器生成部12)は、部分集合が生成される毎に、属性からクラスを判定するルールである分類器(例えば、決定木)を、部分集合に基づいて生成する。

【0122】

10

20

30

40

50

クラス判定手段 8 3 (例えば、自己評価部 1 3) は、分類器が生成される毎に、分類器を用いて、部分集合に属する個々のデータのクラスを判定する。

【0 1 2 3】

類似度算出手段 8 4 (例えば、類似度算出部 1 4) は、データの集合の部分集合が生成され、クラス判定手段がその部分集合に属する個々のデータのクラスを判定したときに、同一のクラスと判定されたデータ同士の類似度に値を加算する。

【0 1 2 4】

本発明によれば、データの重複を許して、データの集合から部分集合を複数回生成し、部分集合毎に、分類器を生成して、部分集合に属するデータのクラスを判定する。そして、この判定結果を用いて類似度を計算するので、データ中に互いに関連のある属性が多く存在してもそれらの関連のある属性の影響を強く受けることなく、データの類似度を求めることができる。

【0 1 2 5】

また、上記の実施形態には、データ同士の類似度に基づいて、データの集合に属するデータを複数のグループに分類するデータグループ化手段 (例えば、類似度クラスタリング部 2 1) を備える構成が開示されている。そのような構成によれば、データの集合に含まれる個別データを分類することができる。

【0 1 2 6】

また、上記の実施形態には、データグループ化手段が、データの集合に属する個々のデータをそれぞれ別々のグループに分類し、互いに異なる二つのグループに属するデータ同士の類似度を求め、類似度が最大となる二つのグループを併合することを繰り返し、グループの総数を目標数まで減少させる構成が開示されている。そのような構成によれば、データの集合に含まれる個別データを目標数のグループに分類することができる。

【0 1 2 7】

また、上記の実施形態には、データ集合に属するデータを、特定の属性の属性値に応じて、グループに分類する属性データ分類手段 (属性データ分類部 3 1) と、データグループ化手段によって分類されたデータのグループと、属性データ分類手段によって分類されたデータのグループとの関係に基づいて、類似度算出に対する特定の属性の関連度を計算する関連度計算手段 (関連性算出部 3 2) とを備える構成が開示されている。そのような構成によれば、類似度を用いて分類を行った結果と属性値との関連の度合いを調べることができる。

【0 1 2 8】

また、上記の実施形態には、部分集合生成手段 8 1 が、データの集合からデータをランダムサンプリングすることによって、集合の部分集合を生成する構成が開示されている。

【0 1 2 9】

また、上記の実施形態には、類似度算出手段 8 4 が、クラス判定手段 8 3 によって特定のクラスと判定されたデータ同士の類似度に対してのみ値を加算する構成が開示されている。そのような構成によれば、特定のクラスのデータの持つ特徴の違いを重視して類似度を算出できる。

【0 1 3 0】

また、上記の実施形態には、分類器生成手段 8 2 が、部分集合に属するデータのうち所与のクラスが特定のクラスであるデータを加重して分類器を生成する構成が開示されている。そのような構成によれば、得られた類似度に基づいてデータを分類するときに、特定のクラス以外のクラスのデータが持つ特徴の違いを重視してグループに分類することができる。

【0 1 3 1】

また、上記の実施形態には、部分集合生成手段 8 1 が、少なくとも顧客の特徴または販売条件を属性とし顧客の行動をクラスとするデータの集合から、その集合の部分集合を生成する構成が開示されている。

【産業上の利用可能性】

10

20

30

40

50

## 【 0 1 3 2 】

本発明は、データの集合に属する各データ間の類似度を求めるデータ類似度計算システムや、データ間の類似度を計算してデータを分類する分類システムに好適に適用される。

## 【図面の簡単な説明】

## 【 0 1 3 3 】

【図 1】本発明の第 1 の実施形態の例を示すブロック図である。

【図 2】データの例を示す説明図である。

【図 3】第 1 の実施形態のデータ類似度計算システムの処理経過の例を示すフローチャートである。

【図 4】分類器の例を示す説明図である。

10

【図 5】分類器から判定されたクラスの例を示す説明図である。

【図 6】本発明の第 2 の実施形態の例を示すブロック図である。

【図 7】最短距離法によって個別データをクラスタリングする処理経過の例を示すフローチャートである。

【図 8】本発明の第 3 の実施形態の例を示すブロック図である。

【図 9】属性の関連度を求める動作の例を示すフローチャートである。

【図 10】 $\chi^2$  分布表の例を示す説明図である。

【図 11】本発明の概要を示すブロック図である。

【図 12】データの例を示す説明図である。

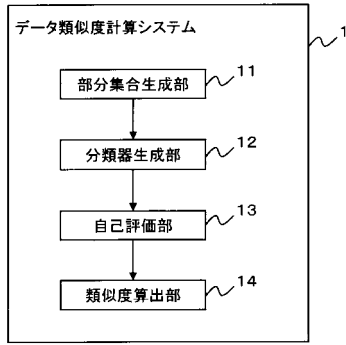
## 【符号の説明】

20

## 【 0 1 3 4 】

- 1 1 部分集合生成部
- 1 2 分類器生成部
- 1 3 自己評価部
- 1 4 類似度算出部
- 2 1 類似度クラスタリング部
- 3 1 属性データ分類部
- 3 2 関連性算出部

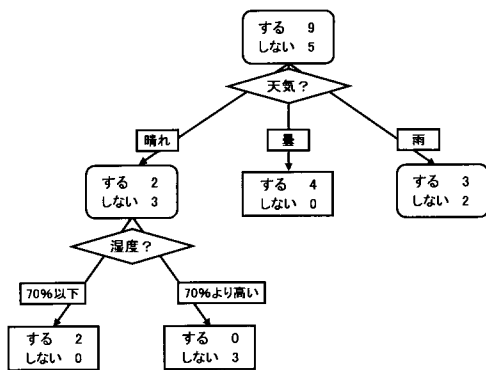
【図1】



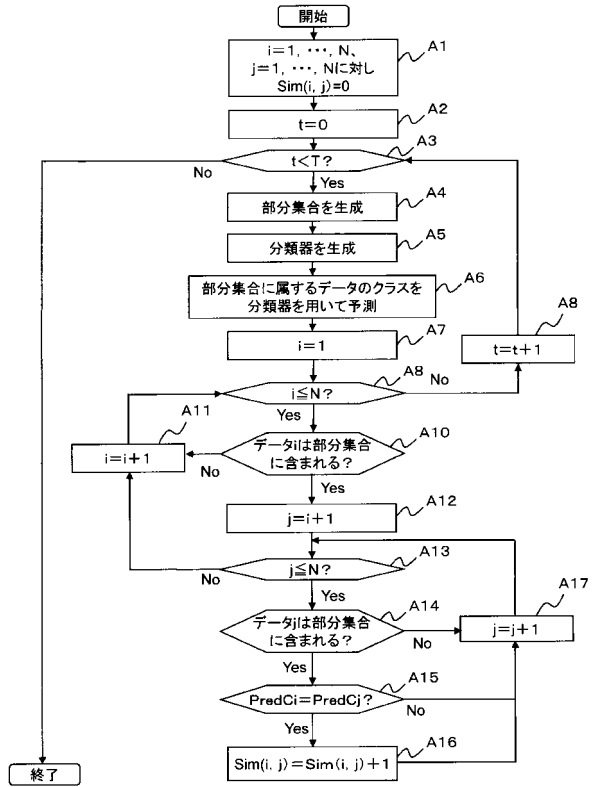
【図2】

属性				クラス
天気	気温(度)	湿度(%)	風が強いかな	購入するか
晴れ	29	85	強くない	しない
晴れ	27	90	強い	しない
曇	28	78	強くない	する
雨	21	96	強くない	する
雨	20	80	強くない	する
雨	18	70	強い	しない
曇	18	65	強い	する
晴	22	95	強くない	しない
晴	21	70	強くない	する
雨	24	80	強くない	する
晴	24	70	強い	する
曇	22	90	強い	する
曇	27	75	強くない	する
雨	22	80	強い	しない

【図4】



【図3】

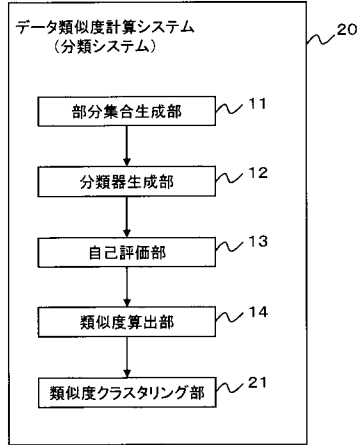


【図5】

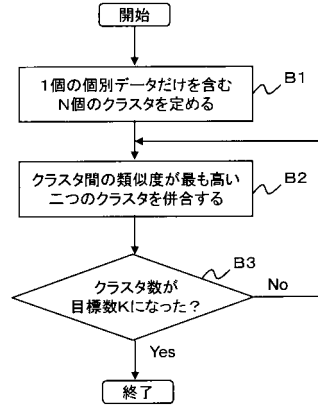
属性				クラス
天気	気温(度)	湿度(%)	風が強いかな	購入するか
晴れ	29	85	強くない	しない
晴れ	27	90	強い	しない
曇	28	78	強くない	する
雨	21	96	強くない	する
雨	20	80	強くない	する
雨	18	70	強い	する
曇	18	65	強い	する
晴	22	95	強くない	しない
晴	21	70	強くない	する
雨	24	80	強くない	する
晴	24	70	強い	する
曇	22	90	強い	する
曇	27	75	強くない	する
雨	22	80	強い	する



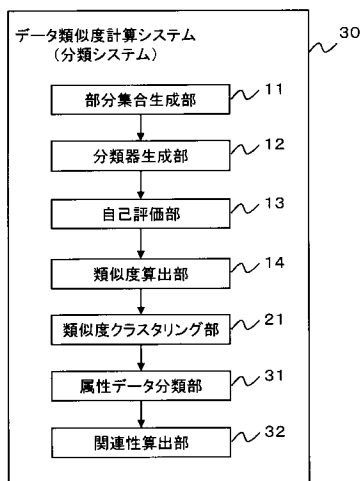
【図6】



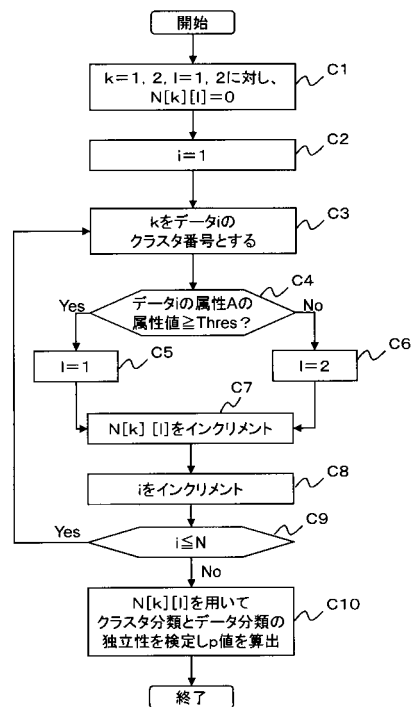
【図7】



【図8】



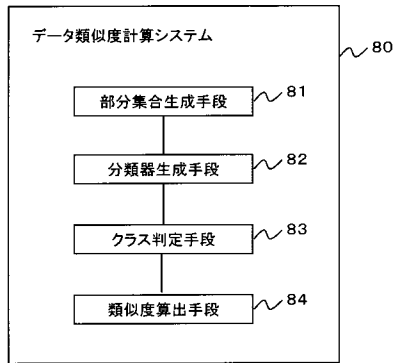
【図9】



【図10】

	a									
y	0.7	0.5	0.3	0.2	0.1	0.05	0.02	0.01	0.001	
1	0.15	0.45	1.07	1.64	2.71	3.84	5.41	6.63	10.83	
2	0.71	1.39	2.41	3.22	4.61	5.99	7.82	9.21	13.82	
3	1.42	2.37	3.66	4.64	6.25	7.81	9.84	11.34	16.27	

【図11】



【図12】

データ番号	属性a	属性b	属性c	属性d	属性e	クラス
1	1	1	1	2	2	A
2	1	1	1	2	2	A
3	1	1	1	2	2	A
4	1	1	1	2	2	A
5	1	2	2	1	1	A
6	2	1	1	2	2	A
7	2	2	2	1	1	A
8	2	2	2	1	1	B
9	2	2	2	1	1	B
10	2	2	2	1	1	B

---

フロントページの続き

(56)参考文献 特開2005-275556(JP,A)

特開2002-149697(JP,A)

特開2005-078240(JP,A)

Leo Breiman, RF/tools A Class of Two-eyed Algorithms, Proceedings of SIAM International Conference on Data Mining, 米国, Society for Industrial and Applied Mathematics, 2003年5月2日, Vol.2003(Keynote4), pp.1-56., [ONLINE]取得日2013.03.07, URL, <http://oz.berkeley.edu/users/breiman/siamtalk2003.pdf>

金明哲、村上征勝、ランダムフォレスト法による文章の書き手の同定, 統計数理, 日本, 統計数理研究所, 2007年, 第55巻、第2号, pp.255-268., [ONLINE]取得日2013/03/07, URL, <http://www.ism.ac.jp/editsec/toukei/pdf/55-2-255.pdf>

Andy Liaw, Matthew Wiener, Classification and Regression by randomForest, R News, オーストリア, the R Foundation for Statistical Computing, communications, 2002年12月, Vol.2/3, pp.18-22., [ONLINE]取得日2013.03.07, URL, [http://cran.r-project.org/doc/Rnews/Rnews\\_2002-3.pdf](http://cran.r-project.org/doc/Rnews/Rnews_2002-3.pdf)

(58)調査した分野(Int.Cl., DB名)

G06F 17/30

JSTPlus(JDreamIII)