

(19) 日本国特許庁(JP)

(12) 特 許 公 報(B2)

(11) 特許番号

特許第4255779号
(P4255779)

(45) 発行日 平成21年4月15日(2009.4.15)

(24) 登録日 平成21年2月6日(2009.2.6)

(51) Int.Cl.		F I	
G 0 6 F 19/00	(2006.01)	G O 6 F 19/00	1 3 0
G 0 6 F 17/15	(2006.01)	G O 6 F 17/15	
G 0 6 F 17/18	(2006.01)	G O 6 F 17/18	Z
G 0 6 Q 50/00	(2006.01)	G O 6 F 17/60	1 0 6

請求項の数 10 (全 31 頁)

(21) 出願番号	特願2003-272648 (P2003-272648)	(73) 特許権者	000005049 シャープ株式会社 大阪府大阪市阿倍野区長池町22番22号
(22) 出願日	平成15年7月10日(2003.7.10)	(74) 代理人	110000338 特許業務法人原謙三国際特許事務所
(65) 公開番号	特開2005-32117 (P2005-32117A)	(74) 代理人	100088281 弁理士 田畑 昌男
(43) 公開日	平成17年2月3日(2005.2.3)	(72) 発明者	竹内 博明 大阪府大阪市阿倍野区長池町22番22号 シャープ株式会社内
審査請求日	平成17年8月10日(2005.8.10)	審査官	小原 正信
		(58) 調査した分野(Int.Cl., DB名)	G 0 6 F 1 9 / 0 0

(54) 【発明の名称】 データ分析装置およびデータ分析方法並びにデータ分析プログラム

(57) 【特許請求の範囲】

【請求項1】

分析対象データ格納部に格納された、複数の入力属性 x_j ($1 \leq j \leq N$ 、 N は入力属性の個数)と、1つの出力属性 y とで構成されるデータの集合である基本データ群 DA を分析対象とし、入力属性と出力属性との因果関係を分析するデータ分析装置であって、

基本データ群 DA に含まれる文字属性のデータを、一義的な変換ルールに従って数値属性のデータに変換することによって、数値属性のデータの集合である数値型基本データ群 $DA0$ を生成する文字 数値データ変換手段と、

数値型基本データ群 $DA0$ を、数値型基本データ群 $DA0$ に含まれる出力属性 y の数値と、出力属性 y の所定閾値との大小関係の比較に基づいて、第1データ群 $DA1$ と、第2データ群 $DA2$ とに分類する分類手段と、

上記複数の入力属性のうちの1つの入力属性 x_j について、該1つの入力属性 x_j のとり得る数値毎に、当該数値以下の数値を持つデータのうち、第1データ群 $DA1$ に属するデータの個数の、第1データ群 $DA1$ に属する全てのデータの個数に対する比率である第1の頻度 ($1 - x_j$ 頻度累積%) を求める演算を行い、かつ、該1つの入力属性 x_j のとり得る数値毎に、当該数値以下の数値を持つデータのうち、第2データ群 $DA2$ に属するデータの個数の、第2データ群 $DA2$ に属する全てのデータの個数に対する比率である第2の頻度 ($2 - x_j$ 頻度累積%) を求める演算を行い、かつ、該1つの入力属性 x_j のとり得る数値毎に、第1の頻度と第2の頻度との差分 (x_j 頻度累積差%) を求める演算を、上記複数の入力属性の各々について行なう第1の評価手段と、

10

20

上記複数の入力属性のうちの1つの入力属性 x_j について、第1の評価手段で該1つの入力属性 x_j のとり得る数値毎に演算された差分 (x_j 頻度累積差%) に基づいて、最大の差分が求められた数値を当該入力属性 x_j の閾値 x_{j_th} として決定することを、上記複数の入力属性の各々について行なう閾値決定手段と、

上記複数の入力属性のうちの1つの入力属性 x_j について、閾値決定手段にて決定された該入力属性 x_j の閾値 x_{j_th} における、第1の頻度 ($1 - x_j$ 頻度累積%) に対する第2の頻度 ($2 - x_j$ 頻度累積%) の比率である第1の比率と、閾値決定手段にて決定された該入力属性 x_j の閾値 x_{j_th} における、($100\% -$ 第1の頻度 ($1 - x_j$ 頻度累積%)) に対する ($100\% -$ 第2の頻度 ($2 - x_j$ 頻度累積%)) の比率である第2の比率とを演算するとともに、第1の比率および第2の比率のうち大きい方の比率を選択することを、上記複数の入力属性の各々について行なう第2の評価手段と、

10

上記第2の評価手段にて入力属性毎に選択された比率のうち、最も大きい比率を持つ入力属性 x_j 、該入力属性 x_j の閾値 x_{j_th} 、および該最も大きい比率が第1の比率および第2の比率の何れであることを示す種別を、入力属性条件を示すデータとして抽出するとともに、当該入力属性条件を分析結果データ格納部に格納する要因抽出手段とを含むことを特徴とするデータ分析装置。

【請求項2】

上記要因抽出手段で抽出された入力属性条件に基づいて、数値型基本データ群 DA0 を、上記入力属性条件を満たす要因データ群と上記入力属性条件を満たさない他データ群とに分割し、分類されたデータ群のうち少なくとも一方を新たな数値型基本データ群 DA0 として分類手段に送る分割手段をさらに含み、

20

分類手段による処理、第1の評価手段による処理、閾値決定手段による処理、第2の評価手段による処理、要因抽出手段による処理、および分割手段による処理からなる一連の処理が繰り返し実行されるようになっていないことを特徴とする請求項1に記載のデータ分析装置。

【請求項3】

上記分割手段は、分類されたデータ群のうち他データ群のみを選択して新たな数値型基本データ群 DA0 として分類手段に送るものであることを特徴とする請求項2に記載のデータ分析装置。

【請求項4】

30

終了条件を満たしているかを判定する終了条件判定手段をさらに含み、上記終了条件判定手段において終了条件を満たしていると判定されると、上記一連の処理の実行を終了するようになっていないことを特徴とする請求項2に記載のデータ分析装置。

【請求項5】

上記終了条件判定手段は、分類手段で分類された第2データ群のデータ数が0であることを終了条件として判定を行なうことを特徴とする請求項4に記載のデータ分析装置。

【請求項6】

予め定められた設定情報に従って、または、使用者からの入力に応じて、出力属性の上記所定閾値を設定する閾値設定手段をさらに含むことを特徴とする請求項1または2に記載のデータ分析装置。

40

【請求項7】

上記入力属性は、製品の製造工程における製造プロセス条件および/またはインライン検査結果であり、

上記出力属性は、製品の品質判定結果であり、

上記第2データ群は、品質判定結果が不良のデータ群であることを特徴とする請求項1または2に記載のデータ分析装置。

【請求項8】

請求項1に記載のデータ分析装置を用いて、分析対象データ格納部に格納された、複数の入力属性 x_j ($1 \leq j \leq N$ 、 N は入力属性の個数) と、1つの出力属性 y とで構成されるデータの集合である基本データ群 DA を分析対象とし、入力属性と出力属性との因果関

50

係を分析するデータ分析方法であって、

上記文字 数値データ変換手段により、基本データ群 D A に含まれる文字属性のデータを、一義的な変換ルールに従って数値属性のデータに変換することによって、数値属性のデータの集合である数値型基本データ群 D A 0 を生成する文字 数値データ変換ステップと、

上記分類手段により、数値型基本データ群 D A 0 を、数値型基本データ群 D A 0 に含まれる出力属性 y の数値と、出力属性 y の所定閾値との大小関係の比較に基づいて、第 1 データ群 D A 1 と、第 2 データ群 D A 2 とに分類する分類ステップと、

上記第 1 の評価手段により、上記複数の入力属性のうちの 1 つの入力属性 x_j について、該 1 つの入力属性 x_j のとり得る数値毎に、当該数値以下の数値を持つデータのうち、第 1 データ群 D A 1 に属するデータの個数の、第 1 データ群 D A 1 に属する全てのデータの個数に対する比率である第 1 の頻度 (1 - x_j 頻度累積 %) を求める演算を行い、かつ、該 1 つの入力属性 x_j のとり得る数値毎に、当該数値以下の数値を持つデータのうち、第 2 データ群 D A 2 に属するデータの個数の、第 2 データ群 D A 2 に属する全てのデータの個数に対する比率である第 2 の頻度 (2 - x_j 頻度累積 %) を求める演算を行い、かつ、該 1 つの入力属性 x_j のとり得る数値毎に、第 1 の頻度と第 2 の頻度との差分 (x_j 頻度累積差 %) を求める演算を、上記複数の入力属性の各々について行なう第 1 の評価ステップと、

上記閾値決定手段により、上記複数の入力属性のうちの 1 つの入力属性 x_j について、第 1 の評価手段で該 1 つの入力属性 x_j のとり得る数値毎に演算された差分 (x_j 頻度累積差 %) に基づいて、最大の差分が求められた数値を当該入力属性 x_j の閾値 x_{j_th} として決定することを、上記複数の入力属性の各々について行なう閾値決定ステップと、

上記第 2 の評価手段により、上記複数の入力属性のうちの 1 つの入力属性 x_j について、閾値決定手段にて決定された該入力属性 x_j の閾値 x_{j_th} における、第 1 の頻度 (1 - x_j 頻度累積 %) に対する第 2 の頻度 (2 - x_j 頻度累積 %) の比率である第 1 の比率と、閾値決定手段にて決定された該入力属性 x_j の閾値 x_{j_th} における、(100 % - 第 1 の頻度 (1 - x_j 頻度累積 %)) に対する (100 % - 第 2 の頻度 (2 - x_j 頻度累積 %)) の比率である第 2 の比率とを演算するとともに、第 1 の比率および第 2 の比率のうち大きい方の比率を選択することを、上記複数の入力属性の各々について行なう第 2 の評価ステップと、

上記要因抽出手段により、上記第 2 の評価手段にて入力属性毎に選択された比率のうち、最も大きい比率を持つ入力属性 x_j 、該入力属性 x_j の閾値 x_{j_th} 、および該最も大きい比率が第 1 の比率および第 2 の比率の何れであることを示す種別を、(補正前請求項 7、段落 0082) 入力属性条件を示すデータとして抽出するとともに、当該入力属性条件を分析結果データ格納部に格納する要因抽出ステップとを含むことを特徴とするデータ分析方法。

【請求項 9】

分析対象データ格納部に格納された、複数の入力属性 x_j (1 ≤ j ≤ N、N は入力属性の個数) と、1 つの出力属性 y とで構成されるデータの集合である基本データ群 D A を分析対象とし、入力属性と出力属性との因果関係を分析するデータ分析装置が備えるコンピュータを機能させるためのデータ分析プログラムであって、

上記データ分析装置は、

基本データ群 D A に含まれる文字属性のデータを、一義的な変換ルールに従って数値属性のデータに変換することによって、数値属性のデータの集合である数値型基本データ群 D A 0 を生成する文字 数値データ変換手段と、

数値型基本データ群 D A 0 を、数値型基本データ群 D A 0 に含まれる出力属性 y の数値と、出力属性 y の所定閾値との大小関係の比較に基づいて、第 1 データ群 D A 1 と、第 2 データ群 D A 2 とに分類する分類手段と、

上記複数の入力属性のうちの 1 つの入力属性 x_j について、該 1 つの入力属性 x_j のとり得る数値毎に、当該数値以下の数値を持つデータのうち、第 1 データ群 D A 1 に属する

10

20

30

40

50

データの個数の、第1データ群DA1に属する全てのデータの個数に対する比率である第1の頻度(1 - x_j 頻度累積%)を求める演算を行い、かつ、該1つの入力属性 x_j のとり得る数値毎に、当該数値以下の数値を持つデータのうち、第2データ群DA2に属するデータの個数の、第2データ群DA2に属する全てのデータの個数に対する比率である第2の頻度(2 - x_j 頻度累積%)を求める演算を行い、かつ、該1つの入力属性 x_j のとり得る数値毎に、第1の頻度と第2の頻度との差分(x_j 頻度累積差%)を求める演算を、上記複数の入力属性の各々について行なう第1の評価手段と、

上記複数の入力属性のうちの1つの入力属性 x_j について、第1の評価手段で該1つの入力属性 x_j のとり得る数値毎に演算された差分(x_j 頻度累積差%)に基づいて、最大の差分が求められた数値を当該入力属性 x_j の閾値 $x_{j_t_h}$ として決定することを、上記複数の入力属性の各々について行なう閾値決定手段と、

上記複数の入力属性のうちの1つの入力属性 x_j について、閾値決定手段にて決定された該入力属性 x_j の閾値 $x_{j_t_h}$ における、第1の頻度(1 - x_j 頻度累積%)に対する第2の頻度(2 - x_j 頻度累積%)の比率である第1の比率と、閾値決定手段にて決定された該入力属性 x_j の閾値 $x_{j_t_h}$ における、(100% - 第1の頻度(1 - x_j 頻度累積%))に対する(100% - 第2の頻度(2 - x_j 頻度累積%))の比率である第2の比率とを演算するとともに、第1の比率および第2の比率のうち大きい方の比率を選択することを、上記複数の入力属性の各々について行なう第2の評価手段と、

上記第2の評価手段にて入力属性毎に選択された比率のうち、最も大きい比率を持つ入力属性 x_j 、該入力属性 x_j の閾値 $x_{j_t_h}$ 、および該最も大きい比率が第1の比率および第2の比率の何れであることを示す種別を、入力属性条件を示すデータとして抽出するとともに、当該入力属性条件を分析結果データ格納部に格納する要因抽出手段とを含み、コンピュータを上記の各手段として機能させるためのデータ分析プログラム。

【請求項10】

請求項9に記載のデータ分析プログラムを記録したコンピュータ読み取り可能な記録媒体。

【発明の詳細な説明】

【技術分野】

【0001】

本発明は、分析対象である出力属性(目的属性)、例えば製造工程で製造される製品の特性等と、出力属性に影響を与える属性である入力属性(説明属性)、例えば製造プロセス条件等との因果関係を分析するデータ分析装置およびデータ分析方法並びにデータ分析プログラムに関する。

【背景技術】

【0002】

出力属性と入力属性との因果関係を分析する有効な手法としては、決定木手法が知られている(特許文献1参照)。この手法では、各入力属性の値で順次切り分けた葉の部分で、出力属性の値がうまくまとまるような木構造を作成する。

【0003】

図10は、特許文献1の従来技術の項(特許文献1の段落[0002]~[0005]および図22参照)に記載されている決定木の1例であり、表1のデータ群を分析対象としている。表1のデータ群は、 x_1 、 x_2 、 x_3 、 x_4 の4つの入力属性の値と、これら入力属性に対する出力属性 y の値とを組とするデータを12個集めた集合である。この手法で作成される決定木(以下、「従来の決定木-1」と呼ぶ事にする)では、図10に示すように、出力属性 y の値 X 、 Y 、 Z が入力属性 x_1 、 x_2 、 x_3 の各値によって、うまく切り分けられている。

【0004】

しかし、図10の従来の決定木-1の作成においては、データを分類する際に、入力属性がとる値の数(属性値の種類数)だけのデータ集合に分類される。例えば、入力属性 x_2 は4種類の値(a 、 b 、 c 、 d)をとるので、入力属性 x_2

10

20

30

40

50

による分類により4つの集合に分類される。そのため、入力属性がとる値の数が増えると、決定木が煩雑になる可能性がある。

【0005】

この課題の解決策として、特許文献1では、各属性において、まとめられる属性値を1つのラベルで表現し、ラベルによりデータ分類する決定木を提案している。

【0006】

図11は、特許文献1の実施例(特許文献1の段落[0010]~[0028]および図13参照)に記載のラベル階層である。この実施例では、例えば、4種の属性値(1, 2, 3, 4)からなる $\times 3$ 属性について、 $\times 3$ 属性値「1」「2」に「2.5以下」というラベルをつけおよび、 $\times 3$ 属性値「3」「4」に「2.5以上」というラベルをつけて階層構造を表現している。このラベル階層構造を用いて作成される決定木(以下、この決定木を従来の決定木-2と呼ぶ事にする)は、図12(特許文献1の段落[0010]~[0028]および図14参照)に示す如くであり、図10に示す従来の決定木-1に比べて、非常に簡潔である。

10

【特許文献1】特開平8-314725号公報(公開日:平成8年(1996)11月29日)

【発明の開示】

【発明が解決しようとする課題】

【0007】

20

上記従来の決定木生成手法をデバイス等の製品の製造工程における製品特性不良の要因分析に応用する場合を題材にして、従来技術の課題を説明する。

【0008】

いま、表1の入力属性 x_1, x_2, x_3, x_4 が製品製造工程における各種のプロセスデータやインライン検査データ、出力属性 y が製造された製品の特性データであり、出力属性 $y = Y$ が製品特性不良に相当するものとする。そして、プロセス技術者が、製品特性不良 $y = Y$ に対し、特許文献1の従来技術に記載された手法で生成された決定木-1(図10)、または特許文献1に記載された手法で生成された従来の決定木-2(図12)を用いて、製品特性不良の要因を調査するものとする。

30

【0009】

このとき、特許文献1の従来技術に記載された手法で生成された決定木-1では、注目すべき $y = Y$ が樹形の中の複数箇所(図10の例では4箇所)に分散しているため煩雑であり、「どの入力属性がどの値の範囲にあるから製品特性が悪いのか?」という製品特性不良の要因をプロセス技術者が判断しにくい。図10の例では、入力属性が4属性だけかつ各属性値の種類も4つだけであるため、何とか、プロセス技術者が製品特性不良の要因を判断することも可能である。しかしながら、実際のデバイス(特に半導体デバイス)のような製品の製造現場では、1工程につき10~100属性程度のプロセスデータやインライン検査データがあり、しかも、その値は多値で非常に広い範囲で分布している。さらに、外乱(入力属性として検出できていない属性)の影響により、各入力属性の値が同じであっても、出力属性の値がばらつく事も多い。これらのような場合に特許文献1の従来技術に記載された手法を用いると、厳密な分析を目指すがあまり、無限数のデータ集合に分類されてしまい、もはや、プロセス技術者が、適正に製品特性不良の要因を特定する事ができなくなる。

40

【0010】

一方、特許文献1に開示された手法により生成される決定木-2(図12)では、ラベル階層による分類がなされているので、決定木が簡潔である。そのため、プロセス技術者が、 $y = Y$ なる製品特性不良の要因を特定しやすい。

【0011】

50

しかし、この図12に示す簡潔な決定木-2を作成するには、図11に示すラベル階層構造を予め定義しておく必要がある。そのため、特許文献1の決定木生成手法は、まとめられる属性値の見当がつかない場合には適用できない。上述したように、実際のデバイスのような製品の製造現場では、1工程につき10~100属性程度の、プロセスデータやインライン検査データがあり、しかも、その値は多値で非常に広い範囲で分布している。さらに、外乱(入力属性として検出できていない属性)の影響により、各入力属性の値が同じであっても、出力属性の値がばらつく事も多い。これらのような状況下で、各入力属性に対し、一つのラベルとしてまとめられる属性値を見出す事は、経験豊富なプロセス技術者であっても、非常に困難である。

10

【0012】

本発明は、上記従来の問題点を鑑みてなされたものであり、その目的は、ラベル階層構造を予め定義する事なく、簡潔な形で、出力属性と入力属性との因果関係を導き出せるデータ分析装置およびデータ分析方法並びにデータ分析プログラムを提供する事にある。

【課題を解決するための手段】

【0013】

本発明に係るデータ分析装置は、上記の課題を解決するために、分析対象データ格納部に格納された、複数の入力属性 x_j ($1 \leq j \leq N$ 、 N は入力属性の個数)と、1つの出力属性 y とで構成されるデータの集合である基本データ群 DA を分析対象とし、入力属性と出力属性との因果関係を分析するデータ分析装置であって、基本データ群 DA に含まれる文字属性のデータを、一義的な変換ルールに従って数値属性のデータに変換することによって、数値属性のデータの集合である数値型基本データ群 $DA0$ を生成する文字数値データ変換手段と、数値型基本データ群 $DA0$ を、数値型基本データ群 $DA0$ に含まれる出力属性 y の数値と、出力属性 y の所定閾値との大小関係の比較に基づいて、第1データ群 $DA1$ と、第2データ群 $DA2$ とに分類する分類手段と、上記複数の入力属性のうちの一つの入力属性 x_j について、該一つの入力属性 x_j のとり得る数値毎に、当該数値以下の数値を持つデータのうち、第1データ群 $DA1$ に属するデータの個数の、第1データ群 $DA1$ に属する全てのデータの個数に対する比率である第1の頻度($1 - x_j$ 頻度累積%)を求める演算を行い、かつ、該一つの入力属性 x_j のとり得る数値毎に、当該数値以下の数値を持つデータのうち、第2データ群 $DA2$ に属するデータの個数の、第2データ群 $DA2$ に属する全てのデータの個数に対する比率である第2の頻度($2 - x_j$ 頻度累積%)を求める演算を行い、かつ、該一つの入力属性 x_j のとり得る数値毎に、第1の頻度と第2の頻度との差分(x_j 頻度累積差%)を求める演算を、上記複数の入力属性の各々について行なう第1の評価手段と、上記複数の入力属性のうちの一つの入力属性 x_j について、第1の評価手段で該一つの入力属性 x_j のとり得る数値毎に演算された差分(x_j 頻度累積差%)に基づいて、最大の差分が求められた数値を当該入力属性 x_j の閾値 x_{j,t_h} として決定することを、上記複数の入力属性の各々について行なう閾値決定手段と、上記複数の入力属性のうちの一つの入力属性 x_j について、閾値決定手段にて決定された該入力属性 x_j の閾値 x_{j,t_h} における、第1の頻度($1 - x_j$ 頻度累積%)に対する第2の頻度($2 - x_j$ 頻度累積%)の比率である第1の比率と、閾値決定手段にて決定された該入力属性 x_j の閾値 x_{j,t_h} における、($100\% -$ 第1の頻度($1 - x_j$ 頻度累積%))に対する($100\% -$ 第2の頻度($2 - x_j$ 頻度累積%))の比率である第2の比率とを演算するとともに、第1の比率および第2の比率のうち大きい方の比率を選択することを、上記複数の入力属性の各々について行なう第2の評価手段と、上記第2の評価手段にて入力属性毎に選択された比率のうち、最も大きい比率を持つ入力属性 x_j 、該入力属性 x_j の閾値 x_{j,t_h} 、および該最も大きい比率が第1の比率および第2の比率の何れであることを示す種別を、入力属性条件を示すデータとして抽出するとともに、当該入力属性条件を分析結果データ格納部に格納する要因抽出手段とを含むことを特徴としている。

20

30

40

50

【0014】

本発明に係るデータ分析方法は、上記の課題を解決するために、前記のデータ分析装置を用いて、分析対象データ格納部に格納された、複数の入力属性 x_j ($1 \leq j \leq N$ 、 N は入力属性の個数)と、1つの出力属性 y とで構成されるデータの集合である基本データ群 DA を分析対象とし、入力属性と出力属性との因果関係を分析するデータ分析方法であって、上記文字 数値データ変換手段により、基本データ群 DA に含まれる文字属性のデータを、一義的な変換ルールに従って数値属性のデータに変換することによって、数値属性のデータの集合である数値型基本データ群 $DA0$ を生成する文字 数値データ変換ステップと、上記分類手段により、数値型基本データ群 $DA0$ を、数値型基本データ群 $DA0$ に含まれる出力属性 y の数値と、出力属性 y の所定閾値との大小関係の比較に基づいて、第1データ群 $DA1$ と、第2データ群 $DA2$ とに分類する分類ステップと、上記第1の評価手段により、上記複数の入力属性のうち1つの入力属性 x_j について、該1つの入力属性 x_j のとり得る数値毎に、当該数値以下の数値を持つデータのうち、第1データ群 $DA1$ に属するデータの個数の、第1データ群 $DA1$ に属する全てのデータの個数に対する比率である第1の頻度 ($1 - x_j$ 頻度累積%)を求める演算を行い、かつ、該1つの入力属性 x_j のとり得る数値毎に、当該数値以下の数値を持つデータのうち、第2データ群 $DA2$ に属するデータの個数の、第2データ群 $DA2$ に属する全てのデータの個数に対する比率である第2の頻度 ($2 - x_j$ 頻度累積%)を求める演算を行い、かつ、該1つの入力属性 x_j のとり得る数値毎に、第1の頻度と第2の頻度との差分 (x_j 頻度累積差%)を求める演算を、上記複数の入力属性の各々について行なう第1の評価ステップと、上記閾値決定手段により、上記複数の入力属性のうち1つの入力属性 x_j について、第1の評価手段で該1つの入力属性 x_j のとり得る数値毎に演算された差分 (x_j 頻度累積差%)に基づいて、最大の差分が求められた数値を当該入力属性 x_j の閾値 x_{j_th} として決定することを、上記複数の入力属性の各々について行なう閾値決定ステップと、上記第2の評価手段により、上記複数の入力属性のうち1つの入力属性 x_j について、閾値決定手段にて決定された該入力属性 x_j の閾値 x_{j_th} における、第1の頻度 ($1 - x_j$ 頻度累積%)に対する第2の頻度 ($2 - x_j$ 頻度累積%)の比率である第1の比率と、閾値決定手段にて決定された該入力属性 x_j の閾値 x_{j_th} における、($100\% -$ 第1の頻度 ($1 - x_j$ 頻度累積%))に対する ($100\% -$ 第2の頻度 ($2 - x_j$ 頻度累積%))の比率である第2の比率とを演算するとともに、第1の比率および第2の比率のうち大きい方の比率を選択することを、上記複数の入力属性の各々について行なう第2の評価ステップと、上記要因抽出手段により、上記第2の評価手段にて入力属性毎に選択された比率のうち、最も大きい比率を持つ入力属性 x_j 、該入力属性 x_j の閾値 x_{j_th} 、および該最も大きい比率が第1の比率および第2の比率の何れであることを示す種別を、入力属性条件を示すデータとして抽出するとともに、当該入力属性条件を分析結果データ格納部に格納する要因抽出ステップとを含むことを特徴としている。

【0015】

本発明に係るデータ分析プログラムは、上記の課題を解決するために、分析対象データ格納部に格納された、複数の入力属性 x_j ($1 \leq j \leq N$ 、 N は入力属性の個数)と、1つの出力属性 y とで構成されるデータの集合である基本データ群 DA を分析対象とし、入力属性と出力属性との因果関係を分析するデータ分析装置が備えるコンピュータを機能させるためのデータ分析プログラムであって、上記データ分析装置は、基本データ群 DA に含まれる文字属性のデータを、一義的な変換ルールに従って数値属性のデータに変換することによって、数値属性のデータの集合である数値型基本データ群 $DA0$ を生成する文字 数値データ変換手段と、数値型基本データ群 $DA0$ を、数値型基本データ群 $DA0$ に含まれる出力属性 y の数値と、出力属性 y の所定閾値との大小関係の比較に基づいて、第1データ群 $DA1$ と、第2データ群 $DA2$ とに分類する分類手段と、上記複数の入力属性のうち1つの入力属性 x_j について、該1つの入力属性 x_j のとり得る数値毎に、当該数値以下の数値を持つデータのうち、第1データ群 $DA1$ に属するデータの個数の、第1データ群 $DA1$ に属する全てのデータの個数に対する比率である第1の頻度 ($1 - x_j$ 頻度累

10

20

30

40

50

積%)を求める演算を行い、かつ、該1つの入力属性 x_j のとり得る数値毎に、当該数値以下の数値を持つデータのうち、第2データ群DA2に属するデータの個数の、第2データ群DA2に属する全てのデータの個数に対する比率である第2の頻度(2- x_j 頻度累積%)を求める演算を行い、かつ、該1つの入力属性 x_j のとり得る数値毎に、第1の頻度と第2の頻度との差分(x_j 頻度累積差%)を求める演算を、上記複数の入力属性の各々について行なう第1の評価手段と、上記複数の入力属性のうち1つの入力属性 x_j について、第1の評価手段で該1つの入力属性 x_j のとり得る数値毎に演算された差分(x_j 頻度累積差%)に基づいて、最大の差分が求められた数値を当該入力属性 x_j の閾値 x_{j_th} として決定することを、上記複数の入力属性の各々について行なう閾値決定手段と、上記複数の入力属性のうち1つの入力属性 x_j について、閾値決定手段にて決定された該入力属性 x_j の閾値 x_{j_th} における、第1の頻度(1- x_j 頻度累積%)に対する第2の頻度(2- x_j 頻度累積%)の比率である第1の比率と、閾値決定手段にて決定された該入力属性 x_j の閾値 x_{j_th} における、(100%-第1の頻度(1- x_j 頻度累積%))に対する(100%-第2の頻度(2- x_j 頻度累積%))の比率である第2の比率とを演算するとともに、第1の比率および第2の比率のうち大きい方の比率を選択することを、上記複数の入力属性の各々について行なう第2の評価手段と、上記第2の評価手段にて入力属性毎に選択された比率のうち、最も大きい比率を持つ入力属性 x_j 、該入力属性 x_j の閾値 x_{j_th} 、および該最も大きい比率が第1の比率および第2の比率の何れであるかを示す種別を、入力属性条件を示すデータとして抽出するとともに、当該入力属性条件を分析結果データ格納部に格納する要因抽出手段とを含み、コンピュータを上記の各手段として機能させるためのデータ分析プログラムであることを特徴としている。

【0016】

本発明に係るコンピュータ読み取り可能な記録媒体は、上記の課題を解決するために、上記のデータ分析プログラムを記録したものであることを特徴としている。

【0017】

上記装置、方法、プログラム、あるいは記録媒体によれば、ラベル階層構造を予め定義する事なく、簡潔な形で、第2データ群に対応する出力属性条件(結果)の要因を抽出できる。それゆえ、例えば第2データ群が悪い結果(例えば不良品の発生)に対応するデータ群であれば、その悪い結果の要因をユーザが容易に把握できる。逆に、第2データ群が良い結果(例えば優れた特性を持つ製品の発生)に対応するデータ群であれば、その良い結果の要因をユーザが容易に把握できる。

【0018】

本発明に係るデータ分析方法は、上記要因抽出手段で抽出された入力属性条件に基づいて、数値型基本データ群DA0を、上記入力属性条件を満たす要因データ群と上記入力属性条件を満たさない他データ群とに分割し、分類されたデータ群のうち少なくとも一方を新たな数値型基本データ群DA0として分類手段に送る分割手段をさらに含み、分類手段による処理、第1の評価手段による処理、閾値決定手段による処理、第2の評価手段による処理、要因抽出手段による処理、および分割手段による処理からなる一連の処理が繰り返し実行されるようになっていくことがより好ましい。

【0019】

上記構成によれば、複数の要因を節点として木構造を作成できる。それゆえ、単独の相関ルールでは表現し難い複数の要因の絡み合った分析対象であっても、十分高い精度で要因を究明できる。

【0020】

本発明に係るデータ分析装置は、終了条件を満たしているかを判定する終了条件判定手段をさらに含み、上記終了条件判定手段において終了条件を満たしていると判定されると、上記一連の処理の実行を終了するようになっていくことがよ

10

20

30

40

50

り好ましい。これにより、必要以上の無駄な処理が行われることを回避できる。

【0021】

上記第1の評価手段は、各入力属性の全ての数値について、第1データ群中における入力属性がその数値以下であるデータの割合を第1の頻度として演算すると共に、第2データ群中における入力属性がその数値以下であるデータの割合を第2の頻度として演算する頻度演算手段と、各入力属性の全ての数値について、第1の頻度と第2の頻度との差分を演算する差分演算手段とを含むことがより好ましい。これにより、閾値評価指標を容易に演算することができる。

【0022】

上記第2の評価手段は、第1のルール評価値として、第1データ群中における入力属性が閾値以下であるデータの割合に対する、第2データ群中における入力属性が閾値以下であるデータの割合の比率を第1の比率として演算すると共に、第2のルール評価値として、第1データ群中における入力属性が閾値を超えるデータの割合に対する、第2データ群中における入力属性が閾値を超えるデータの割合の比率を第2の比率として演算し、双方の比率のうち大きい方の比率を抽出するものであり、上記要因抽出手段は、上記第2の評価手段で抽出された、各入力属性の比率のうちで、その値が最大となる、入力属性、該入力属性の閾値、および抽出された比率の種別を上記入力属性条件を示すデータとして抽出するものであることがより好ましい。これにより、第1および第2のルール評価値を容易に演算することができる。

【発明の効果】

【0023】

本発明の装置、方法、プログラム、記録媒体によれば、以上のように、ラベル階層構造を予め定義する事なく、「入力属性が閾値以下」あるいは「入力属性が閾値を超える」といった非常に簡潔な形で、問題事象である特定の出力属性条件（問題事象）が発生する要因を導き出すことが可能となる。また、複数の要因を導き出せば、それぞれの要因（入力属性）における「入力属性が閾値以下」あるいは「入力属性が閾値を超える」といった条件の組み合わせによる非常に簡潔な形の決定木として、問題事象に関わる因果関係を導き出せる。

【発明を実施するための最良の形態】

【0024】

本発明の一実施形態を以下に説明する。

【0025】

まず、本実施形態のデータ分析装置を図1に基づいて説明する。

【0026】

図1に示すように、データ分析装置は、文字 - 数値データ変換部1、分析対象データ格納部2、閾値設定部（閾値設定手段）3、データ分類部（分類手段）4、データ列抽出部5、頻度演算部（第1の評価手段、頻度演算手段）6、頻度累積差演算部（第1の評価手段、差分演算手段）7、入力属性閾値決定部（閾値決定手段）8、頻度累積比率演算部（第2の評価手段）16、要因抽出部（要因抽出手段）9、要因未発見データ抽出部（分割手段）10、終了条件判定部（終了条件判定手段）11、入力属性閾値テーブル作成部12、寄与率演算部13、分析結果データ格納部14、および出力部15を備えている。

【0027】

次に、次の表1のデータ群DAを分析対象とする場合を例にとり、本実施形態のデータ分析方法を図2に基づいて説明する。表1のデータ群DAは、ハードディスク等の格納部2に格納されている。

【0028】

10

20

30

40

【表 1】

データ群DA

id	x1	x2	x3	x4	y
1	A	a		1	20 X
2	C	a		2	20 X
3	B	b		1	30 X
4	D	b		1	40 X
5	C	a		3	10 Y
6	D	b		3	30 Y
7	A	c		1	10 Y
8	D	d		4	20 Y
9	A	a		4	40 Z
10	B	a		3	30 Z
11	A	b		3	10 Z
12	A	b		4	20 Z

10

【 0 0 2 9 】

表 1 のデータ群 D A は、1 ~ 1 2 の i d (識別子) を持つ 1 2 個のデータから構成されている。表 1 において、x 1 , x 2 , x 3 , x 4 は入力属性である。入力属性 x 1 は 4 つの文字 A , B , C , D のいずれかをとる文字属性である。入力属性 x 2 は 4 つの文字 a , b , c , d のいずれかをとる文字属性である。入力属性 x 3 は 4 つの離散値 1 , 2 , 3 , 4 のいずれかをとる離散属性である。入力属性 x 4 は 4 つの離散値 1 0 , 2 0 , 3 0 , 4 0 のいずれかをとる離散属性である。なお、入力属性は、連続した数値をとる連続属性でもよい。

20

【 0 0 3 0 】

また、表 1 において、y は出力属性である。出力属性は、文字属性であってもよく、離散属性でもよく、また連続属性でもよいが、ここでは、3 つの文字 X , Y , Z のいずれかをとる文字属性である。

【 0 0 3 1 】

本実施形態のデータ分析方法では、y = Y なる場合を問題事象として、出力属性 y が Y となる要因を分析する。

【 0 0 3 2 】

なお、分析対象データの例としては、例えば、入力属性が、製品の製造工程における製造プロセス条件および / またはインライン検査結果 (製造ライン途中での検査結果) 、出力属性が製品の品質判定結果、y = Y なる問題事象が品質判定結果の不良であるデータが挙げられる。この場合、本実施形態のデータ分析方法により入力属性と出力属性との因果関係を分析し、y = Y なる問題事象の要因を導き出すことで、デバイス特性不良等の不良品の発生を解消する対策を容易に図ることが可能となる。したがって、歩留まりの向上等のような製造プロセスの改善を容易に図ることが可能となる。

30

【 0 0 3 3 】

分析対象データのより具体的な例としては、例えば、入力属性 x 1 , x 2 , x 3 , x 4 が、プラズマ C V D プロセスの、ガス流量、ガス圧力、投入電力、成膜時間などのプロセスデータで、出力属性 y が、プラズマ C V D プロセスで形成される薄膜の膜厚であるようなデータが挙げられる。また、これら入力属性および出力属性の値は、連続属性でも離散属性でも文字属性でもよい。文字属性の場合には、例えば、出力属性が膜厚の例で、' 大 ' 、 ' 中 ' 、 ' 小 ' といった具合に表現される。

40

[ステップ 0]

まず、文字 - 数値データ変換部 1 が、ハードディスク等の分析対象データ格納部 2 に格納された表 1 のデータ群 D A における文字属性を下記の変換ルールに従って数値属性 (数値データ) に変換する (S 0) 。これにより、各データは、数値データに変換される。そして、文字 - 数値データ変換部 1 は、変換されたデー

50

タ群をデータ分類部 4 に送る。

(x 1) A 1、 B 2、 C 3、 D 4

(x 2) a 1、 b 2、 c 3、 d 4

(x 3) 変換せず

(x 4) 変換せず

(y) X 1、 Y 2、 Z 3

この変換ルールは、可能な限り、変換後の入力属性の数値が大きいほど出力属性の数値が大きくなるようにあるいはその逆順となるように設定されることが好ましい。なお、変換ルールは、一義性さえあればよく、上記の例に限られない。

【 0 0 3 4 】

上記変換ルールにて数値データに変換されたデータ群 D A 0 は、表 2 に示す通りである。

【 0 0 3 5 】

【表 2】

基本データ群DA0

id	x1	x2	x3	x4	y
1	1	1	1	1	20
2	3	1	2	20	1
3	2	2	1	30	1
4	4	2	1	40	1
5	3	1	3	10	2
6	4	2	3	30	2
7	1	3	1	10	2
8	4	4	4	20	2
9	1	1	4	40	3
10	2	1	3	30	3
11	1	2	3	10	3
12	1	2	4	20	3

【 0 0 3 6 】

この変換により、得られたデータ群 D A 0 は、離散値をとる複数の入力属性 (説明属性) と出力属性 (目的属性) とで構成されるデータの集合となる。以下、データ群 D A 0 を基本データ群と呼ぶ事にする。

[ステップ 1]

閾値設定部 3 は、予め定められた設定情報に従って、あるいは使用者が図示しないキーボードやマウス等の入力部から問題事象の属性値 $y = Y$ を入力したことに応答して、データ群 D A の $y = Y$ なる問題事象に対応する基本データ群 D A 0 の出力属性 y の閾値 (出力属性閾値) y_{th} を設定し、データ分類部 4 に出力する (S 1)。この例においては、データ群 D A の $y = Y$ なる問題事象に対応する基本データ群 D A 0 の出力属性 y の閾値は、 $y_{th} = 2$ である。

[ステップ 2]

次に、データ分類部 4 が、基本データ群 D A 0 の出力属性 y の値と、閾値設定部 3 から出力された出力属性閾値 y_{th} との比較論理 (1) (2) に基づいて、基本データ群 D A 0 を、第 1 データ群 D A 1 と第 2 データ群 D A 2 とに 2 分化 (分類) する (S 2)。

【 0 0 3 7 】

(1) $y > y_{th}$ または $y < y_{th}$ D A 1

(2) $y = y_{th}$ D A 2

言い換えると、データ分類部 4 は、基本データ群 D A 0 を、出力属性が出力属性閾値 y_{th} と一致しない (すなわち 1 または 3 である) 第 1 データ群 D A 1 と、出力属性が出力属性閾値 $y_{th} (= 2)$ と一致する第 2 データ群 D A 2 とに分類する。第 2 データ群 D A 2 は問題事象 (例えば、デバイス特性不良など) のデータ群である。すなわち、第 2 データ群 D A 2 は出力属性 y が問題事象を表す属性値 (

10

20

30

40

50

2)であるデータ群であり、第1データ群DA1は出力属性yが問題事象を表していない属性値(1または3)であるデータ群である。

【0038】

第1データ群DA1を表3に、第2データ群DA2を表4に示す。

【0039】

【表3】

第1データ群DA1

id	x1	x2	x3	x4	y	
1	1	1	1	1	20	1
2	3	1	2	2	20	1
3	2	2	1	3	30	1
4	4	2	1	4	40	1
9	1	1	4	4	40	3
10	2	1	3	3	30	3
11	1	2	3	1	10	3
12	1	2	4	4	20	3

10

【0040】

【表4】

第2データ群DA2

id	x1	x2	x3	x4	y
5	3	1	3	10	2
6	4	2	3	30	2
7	1	3	1	10	2
8	4	4	4	20	2

20

【0041】

なお、以下では、適宜、第1データ群DA1を良品(OK品)データ群、第2データ群DA2を不良品(NG品)データ群と呼ぶ事にする。

[ステップ3]

次に、データ列抽出部5が、良品データ群DA1(表3)から、入力属性x_j(1 ≤ j ≤ 4)の各々のデータ列を抽出する(S3)。このデータ列を1-x_jデータ群と呼ぶ事にする。

30

【0042】

同様に、データ列抽出部5は、不良品データ群DA2(表4)からも、入力属性x_j(1 ≤ j ≤ 4)の各々のデータ列を抽出する(S3)。このデータ列を2-x_jデータ群と呼ぶ事にする。

【0043】

1-x_jデータ群を表5~8に、2-x_jデータ群を表9~12に示す。

【0044】

【表5】

1-x₁データ群

id	x1
1	1
2	3
3	2
4	4
9	1
10	2
11	1
12	1

40

【0045】

【表 6】

1-x2データ群

id	x2
1	1
2	1
3	2
4	2
9	1
10	1
11	2
12	2

10

【 0 0 4 6 】

【表 7】

1-x3データ群

id	x3
1	1
2	2
3	1
4	1
9	4
10	3
11	3
12	4

20

【 0 0 4 7 】

【表 8】

1-x4データ群

id	x4
1	20
2	20
3	30
4	40
9	40
10	30
11	10
12	20

30

【 0 0 4 8 】

【表 9】

2-x1データ群

id	x1
5	3
6	4
7	1
8	4

40

【 0 0 4 9 】

【表 10】

2-x2データ群

id	x2
5	1
6	2
7	3
8	4

【 0 0 5 0 】

【表 1 1】

2-x3データ群

id	x3
5	3
6	3
7	1
8	4

【 0 0 5 1】

【表 1 2】

2-x4データ群

id	x4
5	10
6	30
7	10
8	20

10

【 0 0 5 2】

[ステップ 4]

頻度演算部 6 は、ステップ 3 で良品データ群 D A 1 から抽出された 1 - x j データ群の各々、およびステップ 3 で不良品データ群 D A 2 から抽出された 2 - x j データ群の各々を、入力属性 x j の値で昇順に並べ替える。そして、入力属性 x j の個々の数値について、第 1 データ群におけるその数値以下のデータ個数の割合を表す 1 - x j 頻度累積%と、第 2 データ群におけるその数値以下のデータ個数の割合を表す 2 - x j 頻度累積%とを計算する (S 4)。

20

【 0 0 5 3】

ここでは、表 5 ~ 8 を入力属性 x j の値で昇順に並べ替えた表 1 3 ~ 1 6 を用い、各行 (i d) のデータについて表中でそのデータの位置以上の位置にあるデータ個数の、第 1 データ群の全データ数 (= 8) に対する割合を 1 - x j 頻度累積%として計算している。同様に、表 9 ~ 1 2 を入力属性 x j の値で昇順に並べ替えた表 1 7 ~ 2 0 を用い、各行 (i d) のデータについて表中でそのデータの位置以上の位置にあるデータ個数の、第 2 データ群の全データ数 (= 4) に対する割合を 2 - x j 頻度累積%として計算している

30

ここで計算した 1 - x j 頻度累積%および 2 - x j 頻度累積%の値を表 1 3 ~ 2 0 に示す。

【 0 0 5 4】

【表 1 3】

1-x1データ群

id	x1	1-x1頻度 累積%
1	1	12.5
9	1	25
11	1	37.5
12	1	50
3	2	62.5
10	2	75
2	3	87.5
4	4	100

40

【 0 0 5 5】

【表 1 4】

1-x2データ群

id	x2	1-x2頻度 累積%
1	1	12.5
2	1	25
9	1	37.5
10	1	50
3	2	62.5
4	2	75
11	2	87.5
12	2	100

10

【 0 0 5 6】

【表 1 5】

1-x3データ群

id	x3	1-x3頻度 累積%
1	1	12.5
3	1	25
4	1	37.5
2	2	50
10	3	62.5
11	3	75
9	4	87.5
12	4	100

20

【 0 0 5 7】

【表 1 6】

1-x4データ群

id	x4	1-x4頻度 累積%
11	10	12.5
1	20	25
2	20	37.5
12	20	50
3	30	62.5
10	30	75
4	40	87.5
9	40	100

30

【 0 0 5 8】

【表 1 7】

2-x1データ群

id	x1	2-x1頻度 累積%
7	1	25
5	3	50
6	4	75
8	4	100

40

【 0 0 5 9】

【表 1 8】

2-x2データ群

id	x2	2-x2頻度 累積%
5	1	25
6	2	50
7	3	75
8	4	100

【 0 0 6 0 】

【表 1 9】

2-x3データ群

id	x3	2-x3頻度 累積%
7	1	25
5	3	50
6	3	75
8	4	100

10

【 0 0 6 1 】

【表 2 0】

2-x4データ群

id	x4	2-x4頻度 累積%
5	10	25
7	10	50
8	20	75
6	30	100

20

【 0 0 6 2 】

なお、上述したステップ 3・4 では、データ列を抽出し、並び替えを行った後に、1-xj 頻度累積%および 2-xj 頻度累積%を計算していたが、データ列の抽出や並び替えを行うことなく直接的に 1-xj 頻度累積%および 2-xj 頻度累積%を計算してもかまわない。

【 0 0 6 3 】

さらに、頻度演算部 6 は、1-xj 頻度累積%が計算された良品データ群である 1-xj データ群のテーブルと、2-xj 頻度累積%が計算された不良品データ群である 2-xj データ群のテーブルとを結合する。具体的には、入力属性 x1 について、表 1 3 と表 1 7 とを結合して表 2 1 の x1 頻度累積テーブルを、入力属性 x2 について、表 1 4 と表 1 8 とを結合して表 2 2 の x2 頻度累積テーブルを、入力属性 x3 について、表 1 5 と表 1 9 とを結合して表 2 3 の x3 頻度累積テーブルを、入力属性 x4 について、表 1 6 と表 2 0 とを結合して表 2 4 の x4 頻度累積テーブルを、それぞれ作成する。

30

【 0 0 6 4 】

【表 2 1】

x1 頻度累積テーブル

id	x1	1-x1 頻度 累積%	2-x1 頻度 累積%
1	1	12.5	
9	1	25	
11	1	37.5	
12	1	50	
3	2	62.5	
10	2	75	
2	3	87.5	
4	4	100	
7	1		25
5	3		50
6	4		75
8	4		100

10

【 0 0 6 5】

【表 2 2】

x2 頻度累積テーブル

id	x2	1-x2 頻度 累積%	2-x2 頻度 累積%
1	1	12.5	
2	1	25	
9	1	37.5	
10	1	50	
3	2	62.5	
4	2	75	
11	2	87.5	
12	2	100	
5	1		25
6	2		50
7	3		75
8	4		100

20

30

【 0 0 6 6】

【表 2 3】

x3 頻度累積テーブル

id	x3	1-x3 頻度 累積%	2-x3 頻度 累積%
1	1	12.5	
3	1	25	
4	1	37.5	
2	2	50	
10	3	62.5	
11	3	75	
9	4	87.5	
12	4	100	
7	1		25
5	3		50
6	3		75
8	4		100

40

50

【 0 0 6 7 】

【 表 2 4 】

x4頻度累積テーブル

id	x4	1-x4頻度 累積%	2-x4頻度 累積%
11	10	12.5	
1	20	25	
2	20	37.5	
12	20	50	
3	30	62.5	
10	30	75	
4	40	87.5	
9	40	100	
5	10		25
7	10		50
8	20		75
6	30		100

10

【 0 0 6 8 】

さらに、頻度演算部6は、表21～24の各々の頻度累積テーブルを、入力属性xjの値で昇順に並べ替える。このとき、1-xj頻度累積%および2-xj頻度累積%の空欄には、その直前の値を代入する。また、入力属性xjにおいて同じ値が続いている場合には、上記並べ替えられた最終のデータのみを採用する。こうして、頻度演算部6にて、入力属性xjの各値に対して、良品データ群である第1データ群におけるその数値以下のデータ個数の割合を表す1-xj頻度累積%(A;第1の頻度)と、不良品データ群である第2データ群におけるその数値以下のデータ個数の割合を表す2-xj頻度累積%(B;第2の頻度)との双方が算出される(S4)。

20

【ステップ5】

次に、頻度累積差演算部7が、入力属性xjの各値に対して、良品の1-xj頻度累積(A)と、不良品の2-xj頻度累積(B)の差分(=|A-B|)を計算する(S5)。この差分値を、xj頻度累積差(=|A-B|)と呼ぶ。xj頻度累積差の計算結果を表25～表28に示す。

30

【 0 0 6 9 】

【 表 2 5 】

x1頻度累積テーブル

x1	1-x1頻度 累積% (A)	2-x1頻度 累積% (B)	x1頻度累 積差% (A-B)
	1	50	25
x1-th→	2	75	50
	3	87.5	37.5
	4	100	0

40

【 0 0 7 0 】

【表 2 6】

x2頻度累積テーブル

x2	1-x2頻度 累積% (A)	2-x2頻度 累積% (B)	x2頻度累 積差% (A-B)
1	50	25	25
2	100	50	50
3	100	75	25
4	100	100	0

x2-th→

10

【 0 0 7 1】

【表 2 7】

x3頻度累積テーブル

x3	1-x3頻度 累積% (A)	2-x3頻度 累積% (B)	x3頻度累 積差% (A-B)
1	37.5	25	12.5
2	50	25	25
3	75	75	0
4	100	100	0

x3-th→

20

【 0 0 7 2】

【表 2 8】

x4頻度累積テーブル

x4	1-x4頻度 累積% (A)	2-x4頻度 累積% (B)	x4頻度累 積差% (A-B)
1	12.5	50	37.5
2	50	75	25
3	75	100	25
4	100	100	0

x4-th→

30

【 0 0 7 3】

入力属性 x_j と、良品の $1 - x_j$ 頻度累積 (A)、不良品の $2 - x_j$ 頻度累積 (B)、 x_j 頻度累積差 $|A - B|$ との関係を図 3 ~ 図 6 に示す。

40

【 0 0 7 4】

各数値に対する x_j 頻度累積差 $|A - B|$ は、入力属性 x_j がその数値以下の範囲と、入力属性 x_j がその数値を超える範囲との 2 分化によって、良品の第 1 データ群 DA1 と不良品の第 2 データ群 DA2 とがうまく切り分けられているかを表す指標である。言い換えると、 x_j 頻度累積差 $|A - B|$ は、入力属性がその数値以下であるデータが第 1 データ群および第 2 データ群のうち一方に偏っている度合いを表す閾値評価指標である。

【 0 0 7 5】

なお、ここでは、閾値評価指標として x_j 頻度累積差 $|A - B|$ を演算してい

50

るが、各数値に対する閾値評価指標として、データの偏りの度合いを評価する指標、例えば、情報利得（ゲイン）、情報利得比、Giniインデックス、平均自乗誤差等を用いてもよい。

【ステップ6】

入力属性閾値決定部8が、各入力属性 x_j について、 x_j の個々の値の中で、 x_j 頻度累積差 $|A - B|$ の値が最大となる時の入力属性 x_j の値を抽出する（S6）。この値を、入力属性閾値 x_{j-th} と呼ぶ事にする。

【0076】

入力属性閾値 x_{j-th} は、図3～図6を参照して分かるように、 $x_j = x_{j-th}$ の範囲と、 $x_j > x_{j-th}$ の範囲との2分化によって、良品の第1データ群DA1と、不良品の第2データ群DA2との切分けが最も容易となる入力属性 x_j の値を示している。

10

【0077】

なお、ここでは、複数の入力属性について第3ステップ～第6ステップの処理を一括して行っているが、 j の値を1からNまで順次増加させて第3ステップ～第6ステップの処理を繰り返してもよい。

【ステップ7】

次に、頻度累積比率演算部16が、 $x_j = x_{j-th}$ において、良品の1 - x_j 頻度累積（A）に対する、不良品の2 - x_j 頻度累積（B）の比率を計算する。この比率を、2 - x_{j-th} 下比率（ $= B / A$ ）と呼ぶ事にする。また、100から良品の1 - x_j 頻度累積（A）を引いた値（ $= 100 - A$ ）に対する、100から不良品の2 - x_j 頻度累積（B）を引いた値（ $= 100 - B$ ）の比率を計算する。この比率を、2 - x_{j-th} 上比率（ $= (100 - B) / (100 - A)$ ）と呼ぶ事にする。そして、双方の比率のうち大きい方の値を表す、2 - x_{j-th} 比率を抽出する。

20

【0078】

ここで、2 - x_{j-th} 下比率は、「 $x_j = x_{j-th}$ 」という入力属性条件により、良品の第1データ群と分離して不良品の第2データ群を検出できる割合を表している。また、2 - x_{j-th} 上比率は、「 $x_j > x_{j-th}$ 」という入力属性条件により、良品の第1データ群と分離して不良品の第2データ群を検出できる割合を表している。

30

【0079】

言い換えると、2 - x_{j-th} 下比率は、「入力属性 x_j が入力属性閾値 x_{j-th} 以下であれば第2データ群に含まれるデータである」という第1の相関ルールの確からしさを表す評価値（第1のルール評価値）を表している。また、2 - x_{j-th} 上比率は、「入力属性 x_j が入力属性閾値 x_{j-th} を超えていれば第2データ群に含まれるデータである」という第2の相関ルールの確からしさを表す評価値（第2のルール評価値）を表している。

【0080】

各入力属性 x_j に対して抽出された入力属性閾値 x_{j-th} 、 $x_j = x_{j-th}$ における、良品の1 - x_j 頻度累積（A）、不良品の2 - x_j 頻度累積（B）、 x_j 頻度累積差 $|A - B|$ 、2 - x_{j-th} 下比率 B / A 、2 - x_{j-th} 上比率 $(100 - B) / (100 - A)$ 、2 - x_{j-th} 比率の各値を表29に示す。

40

【0081】

【表 29】

入力属性 x_j	入力属性 閾値 x_{j-th}	1- x_j 頻度 累積% (A)	2- x_j 頻度 累積% (B)	x_j 頻度累 積差% (A-B)	2- x_j th下 比率 (B/A)	2- x_j th上 比率 (((100-B)/ (100-A))	2- x_j th比 率
1	2	75	25	50	0.333333	3	3
2	2	100	50	50	0.5	∞	∞
3	2	50	25	25	0.5	1.5	1.5
4	1	12.5	50	37.5	4	0.571429	4

【0082】

10

[ステップ8]

要因抽出部9が、 $x_1 \sim x_4$ の入力属性のうち、上記ステップ7の2 - x_j t h 比率が最大となる入力属性を抽出する。これにより、2 - x_j t h 比率が最大となる入力属性と、その閾値、採用した比率の種別(上、下)が第2データ群に対応する出力属性条件の要因(入力属性条件)を示すデータとして抽出されることになる。これは、全ての入力属性に関する前記相関ルールのうちで最も高い2 - x_j t h 下比率または2 - x_j t h 上比率を持つ相関ルールの入力属性条件を示すデータを抽出することに相当する。

【0083】

20

なお、ここでは、最大のルール評価値を持つ相関ルールの入力属性を抽出するための指標として2 - x_j t h 比率を演算しているが、最大のルール評価値を持つ相関ルールの入力属性を抽出するための指標として、他の評価指標、例えば、支持率(サポート)、確信度(コンフィデンス)、情報利得(ゲイン)、情報利得比、Giniインデックス、平均自乗誤差等を用いてもよい。

【0084】

表29を参照して、入力属性 $x_2 = x_2 - th = 2$ のとき、2 - x_2 t h 比率 = 2 - x_2 t h 上比率 = となっている。これは、入力属性条件「 $x_2 > 2$ 」にて、良品の第1データ群DA1と完全に分離して、不良品の第2データ群DA2を検出できる事を示しており、この事は、図4を参照すると、より理解しやすい。

30

【0085】

上記抽出された、入力属性(= x_2)、該入力属性の値を表す入力属性閾値(=2)、および採用した比率の種別(=上)のデータを分析結果データ格納部14に保存する。

【0086】

以上のようにして、問題事象(不良品の第2データ群DA2)の一要因として、「 $x_2 > 2$ 」という入力属性条件が抽出された。

[ステップ9]

上記ステップ8にて、問題事象(不良品の第2データ群DA2)の一要因として、「 $x_2 > 2$ 」という入力属性条件が抽出されたので、次に、別の要因を調査する。このため、要因未発見データ抽出部10が、基本データ群DA0(表2)を入力属性条件「 $x_2 > 2$ 」を満たすデータ群(要因データ群)と、基本データ群DA0(表2)の中で問題事象の要因をまだ発見できていないデータ群(他データ群)、すなわち入力属性条件「 $x_2 \leq 2$ 」を満たす(入力属性条件「 $x_2 > 2$ 」を満たさない)データ群とに分割し、問題事象の要因をまだ発見できていないデータ群を抽出する(表30)。

40

【0087】

【表 3 0】

基本データ群DA0(2回目)

id	x1	x2	x3	x4	y
1	1	1	1	20	1
2	3	1	2	20	1
3	2	2	1	30	1
4	4	2	1	40	1
5	3	1	3	10	2
6	4	2	3	30	2
9	1	1	4	40	3
10	2	1	3	30	3
11	1	2	3	10	3
12	1	2	4	20	3

10

【 0 0 8 8 】

要因未発見データ抽出部 1 0 は、抽出されたデータ群を次の（新しい）基本データ群 D A 0 としてデータ分類部 4 に送る。

[ステップ 1 0]

そして、ステップ 9 で抽出されたデータ群を次の基本データ群 D A 0 として、終了条件判定部 1 1 で終了条件を満たしていると判定されるまで、上記のステップ 2 ~ ステップ 9 の処理が繰り返される。本実施形態の終了条件判定部 1 1 は、繰り返し処理中の上記ステップ 2 において不良品の第 2 データ群 D A 2 のデータ個数が 0 となった場合を終了条件と判定するようになっている。このように不良品の第 2 データ群 D A 2 のデータ個数が 0 となるまで繰り返し処理を実行することにより、詳細な要因分析結果が得られる。

20

【 0 0 8 9 】

なお、終了条件は、第 2 データ群 D A 2 のデータ個数に基づく他の終了条件、例えば、（ 1 ）繰り返し処理中の上記ステップ 2 において第 2 データ群 D A 2 のデータ個数が所定数以下となった場合、（ 2 ）繰り返し処理中の上記ステップ 2 において第 1 データ群 D A 1 のデータ個数に対する第 2 データ群 D A 2 のデータ個数の割合が所定割合以下となった場合、（ 3 ）繰り返し処理中の上記ステップ 8 において抽出された入力属性条件のルール評価値が所定の値を下回った場合等としてもよい。これらのような終了条件を用いた場合、より簡潔で十分な要因分析結果を得ることができる。さらに、簡潔な要因分析結果を得ることを優先する場合には、終了条件を単に繰り返し処理を所定回数行った場合としたり、終了条件判定部 1 1 を省いて、可能な限り繰り返し処理を行うようにしてもよい。

30

【 0 0 9 0 】

今回の例では、2 回目の繰り返し処理中のステップ 9 で抽出した、要因未発見の、x 1 2 のデータ群に不良品のデータ（第 2 データ群 D A 2 ; y = 2 ）が含まれていなかったため、繰り返し処理は 2 回目で（ 2 回目の要因抽出を行った時点で）終了した。

[ステップ 1 1]

入力属性閾値テーブル作成部 1 2 が、ステップ 1 0 の繰り返し処理毎に抽出された入力属性 x j と、入力属性閾値 x j - t h と、採用された比率の種別とを格納した入力属性閾値テーブルを作成する（表 3 1 ）。

40

【 0 0 9 1 】

【表 3 1】

入力属性閾値テーブル

	入力属性 x _j	入力属性 閾値 x _j -th	種別	要因条件
1回目	x2	2	上	x2>2
2回目	x1	2	上	x1>2

50

【 0 0 9 2 】

入力属性閾値テーブル作成部 1 2 では、必要に応じて、入力属性閾値テーブルにおける入力属性閾値 $x_j - t h$ の数値を文字データに変換する。文字データへの変換ルールは、ステップ 0 の変換の逆変換となるルールであり、下記の通りである。

- (x 1) 1 A、 2 B、 3 C、 4 D
- (x 2) 1 a、 2 b、 3 c、 4 d
- (x 3) 変換せず
- (x 4) 変換せず

表 3 1 の入力属性閾値テーブルにおける入力属性閾値 $x_j - t h$ を文字データに変換した入力属性閾値テーブルを表 3 2 に示す。

10

【 0 0 9 3 】

【表 3 2】

入力属性閾値テーブル

	入力属性 x_j	入力属性 閾値 $x_j - t h$	種別	要因条件
1回目	x2	b	上	$x2=c$ or d
2回目	x1	B	上	$x1=C$ or D

【 0 0 9 4 】

この入力属性閾値テーブルは、特許文献 1 に記載の従来決定木 - 2 (図 1 2) において、出力属性 $y = Y (y = 2)$ の切分けに着目した場合の決定木の分類条件に対応する。

20

[ステップ 1 2]

次に、寄与率演算部 1 3 が、表 3 1 の入力属性閾値テーブルから、抽出された入力属性の、問題事象 ($y = 2$: 不良品データ群である、元の第 2 データ群 D A 2) に対する寄与率 (相関ルールの評価指標であるサポートに相当する) を求める。

【 0 0 9 5 】

表 3 3 は、問題事象 (不良品) である元の第 2 データ群 D A 2 (表 4) において、その要因として 1 回目に抽出された「 $x 2 > 2$ 」なる入力属性条件、または、2 回目に抽出された「 $x 1 > 2$ 」なる入力属性条件、に該当するデータに「 * 」を付したものである。

30

【 0 0 9 6 】

【表 3 3】

第 2 データ群 D A 2

id	x1	x2	x3	x4	y
5	*3	1	3	10	2
6	*4	2	3	30	2
7	1	*3	1	10	2
8	*4	*4	4	20	2

40

【 0 0 9 7 】

表 3 3 から、問題事象 (元の第 2 データ群 D A 2) に対する入力属性条件「 $x 1 > 2$ 」、「 $x 2 > 2$ 」の寄与率が表 3 4 に示すように求められる。

【 0 0 9 8 】

【表 3 4】

問題事象に対する要因寄与率

	$x1 > 2$	$x2 > 2$	Total
$x1 > 2$	50%	25%	75%
$x2 > 2$	25%	25%	50%

【 0 0 9 9 】

50

表34において、「 $x_1 > 2$ 」と「 $x_1 > 2$ 」との交差点に示す寄与率、及び「 $x_2 > 2$ 」と「 $x_2 > 2$ 」との交差点に示す寄与率は、それぞれ「 $x_1 > 2$ 」単独要因の寄与率、及び「 $x_2 > 2$ 」単独要因の寄与率を、それぞれ表している。また、「 $x_1 > 2$ 」と「 $x_2 > 2$ 」との交差点に示す寄与率は何れも、「 $x_1 > 2$ 」要因と「 $x_2 > 2$ 」要因との複合要因の寄与率を表している。なお、表34は、図7のようにも表現できる。

【0100】

表34または図7から、問題事象 ($y = 2$) に対し、優先順位 (順位1: x_1 , 順位2: x_2) を付けて対策を施す事ができる。

[ステップ13]

以上でデータ分析を終了し、入力属性閾値テーブル作成部12で作成された入力属性閾値テーブルや、寄与率のデータが、分析結果データとしてハードディスク等の分析結果データ格納部14に格納される。この分析結果データは、適宜、分析結果データ格納部14から表示装置や印刷装置等の出力部15に送られ、表示装置にて決定木やテーブルとして表示したり、印刷装置にて決定木やテーブルとして印刷したりすることができる。

【0101】

本実施形態によれば、特許文献1に記載の、従来の決定木-2 (図12) のように、ラベル階層構造 (図11) を予め定義しなくても、表32 (または表31) の入力属性閾値テーブルに示したような非常に簡潔な形で、問題事象の要因を導き出せる。そして、これを用いて、問題事象に対する各要因 (入力属性) の寄与率を求める事ができる。

【0102】

ここで、表32 (または表31) に示される本実施形態の入力属性閾値テーブルを、決定木の形式で表現すると、図8のように表される。また、従来の決定木-2 (図12) を用いて、図7と同じ形式で、問題事象 $y = Y (= 2)$ に対する各要因の寄与率を表現すると、図9のようになる。

【0103】

本実施形態から導かれる決定木 (図8) と、従来の決定木-2 (図12) とを比較すると、本実施形態の場合には、入力属性 x_3 の寄与が表現されていない。これは、図7と図9とを比較して分かるように、問題事象 $y = Y (y = 2)$ が、入力属性 x_1 および x_3 の、それぞれの単独要因では発生していないからであり、上記の2回目の繰り返し操作中のステップ9において、 $x_1 > 2$ のデータ群に対してステップ10を実行しなかった事に因る。

【0104】

詳細に要因を追求する場合には、入力属性 x_3 の寄与も抽出する必要があるが、問題事象 $y = Y (y = 2)$ を除去する (改善する) 事を目的すれば、入力属性 x_1 のみの抽出であってもこの目的を十分に達成できる。本実施形態では、この点に着目し、問題事象に対して対策すべき主要因を抽出しているため、入力属性 x_3 を抽出していない。詳細な分析を必要とする場合には、上記ステップ9で2分化されたデータ群の双方に対して、ステップ10を実行すればよい。

【0105】

なお、上述した実施形態では、複数の要因を導き出し決定木を生成していたが、単に一つの要因だけを抽出したい場合であれば、ステップ8で終了してもよい。

【0106】

以上で説明したデータ分析方法は、コンピュータが図2のS0~S12 (ステップ0~13) に対応するプロセスを含むデータ分析プログラムを実行することによって実現できる。したがって、図1のデータ分析装置は、データ分析プログラムが、コンピュータを文字-数値データ変換部1、分析対象データ格納部2、

10

20

30

40

50

閾値設定部 3、データ分類部 4、データ列抽出部 5、頻度演算部 6、頻度累積差演算部 7、入力属性閾値決定部 8、頻度累積比率演算部 16、要因抽出部 9、要因未発見データ抽出部 10、終了条件判定部 11、入力属性閾値テーブル作成部 12、および寄与率演算部 13として機能させることにより実現することが可能である。

【0107】

上記プログラムは、コンピュータで読み取り可能な記録媒体に格納してユーザに提供することができる。この記録媒体は、コンピュータ本体に内蔵された内蔵メディアであってもよいし、コンピュータ本体に対して分離可能に構成されたリムーバブル・メディアであってもよい。上記内蔵メディアとしては、ROM；フラッシュメモリ等の書き換え可能な不揮発性メモリ；ハードディスク等が挙げられる。また、上記リムーバブル・メディアとしては、CD-ROM、DVD等の光記録媒体；MO等の光磁気記録媒体；フロッピー（登録商標）ディスク、カセットテープ、リムーバブル・ハードディスク等の磁気記録媒体；メモリカード等のような書き換え可能な不揮発性メモリを内蔵したメディア；ROMカセット等のようなROMを内蔵したメディア等が挙げられる。

10

【0108】

上記プログラムは、CPUのアクセスにより実行される構成であってもよいし、記録媒体に格納されているプログラムを読み出し、読み出したプログラムを内蔵メディアのプログラム記憶領域に転送した後、内蔵メディア上のプログラムがCPUのアクセスにより実行される構成であってもよい。また、上記プログラムは、コンピュータで読み取り可能な記録媒体に格納された状態で販売されるものに限定されるものではなく、インターネット等の通信ネットワークを介してユーザのコンピュータに転送する形式で販売されるものであってもよい。

20

【0109】

なお、本実施形態では、データ分類部 4において出力属性と出力属性閾値との比較により分類を行っていたが、出力属性が文字属性である場合、文字 - 数値データ変換部 1で出力属性を数値属性に変換せず、データ分類部 4において出力属性と要因分析対象となる出力属性（文字；Y）との比較により分類を行うようにしてもよい。

30

【0110】

本実施形態に係るデータ分析方法は、以上のように、N個（Nは2以上の整数）の属性からなるN列の入力属性のデータと、1個の属性からなる1列の出力属性のデータとで構成される基本データ群を分析対象とし、該出力属性と該入力属性との因果関係を分析するデータ分析方法であって、出力属性閾値を設定する第1ステップと、該出力属性の値と該出力属性閾値との比較に基づいて、該基本データ群を、第1データ群と第2データ群とに2分化する第2ステップと、該第1データ群および該第2データ群の各々から、第J入力属性（Jは、1 ≤ J ≤ Nなる関係にある整数）のデータ列を表す1-Jデータ列および2-Jデータ列を、それぞれ抽出する第3ステップと、該1-Jデータ列の該第J入力属性の個々の値に対して、その値以下のデータ個数の割合を表す1-J頻度累積（％）を計算し、該2-Jデータ列の該第J入力属性の個々の値に対して、その値以下のデータ個数の割合を表す2-J頻度累積（％）を計算する第4ステップと、該1-Jデータ列および該2-Jデータ列の双方を含めた、該第J入力属性の全ての値の個々に対して、該1-J頻度累積（％）と該2-J頻度累積（％）との差の絶対値を表す、第J頻度累積差を計算する第5ステップと、第J頻度累積差の値が最大となるときの第J入力属性の値を第J入力属性閾値として抽出する第6ステップと、第J入力属性が第J入力属性閾値であるときにおいて、該1-J頻度累積（％）に対する該2-J頻度累積（％）の比率を表す2-J下比率、および、100から該1-J頻度累積（％）を引いた値に対する、100から該2-J頻度

40

50

累積(%)を引いた値の比率を表す2-J上比率を計算し、双方の比率のうちの大きい方の値を示す、2-J比率を抽出する第7ステップと、Jの値を1からNまで順次増加させて、該第3ステップ~該第7ステップの操作を繰り返し、繰り返し操作中の該第7ステップで抽出された、第1から第Nまでの入力属性の該2-J比率のうち、その値が最大となる入力属性、該入力属性の値を表す入力属性閾値、および採用した比率の種別を抽出し、保存する第8ステップと、該第8ステップで抽出された入力属性に基づいて、該基本データ群を2分化する第9ステップと、該第9ステップで2分化されたデータ群のうちの少なくとも一方を、新たな基本データ群として、所定の終了条件を満たすまで、該第2ステップ~該第9ステップの操作を繰返す第10ステップとを含む。

10

【0111】

上記方法によれば、ラベル階層構造を予め定義しなくても、非常に簡潔な形で問題事象の要因を複数導き出せる。そして、これを用いて、因果関係を表す決定木を作成したり、問題事象(出力属性)に対する各要因(入力属性)の寄与率を求めたりすることができる。

なお、本発明に係るデータ分析装置は、上記の課題を解決するために、複数の入力属性と、出力属性とで構成されるデータの集合である基本データ群を分析対象とし、入力属性と出力属性との因果関係を分析し、因果関係を示す情報を抽出するデータ分析装置であって、基本データ群を出力属性に依って第1データ群と第2データ群とに分類する分類手段と、各入力属性の全ての数値について、入力属性がその数値以下であるデータが第1データ群および第2データ群のうち一方に偏っている度合いを表す閾値評価指標を演算する第1の評価手段と、第1の評価手段で演算された閾値評価指標に基づいて、各入力属性について最大の閾値評価指標を持つ数値を各入力属性の閾値として決定する閾値決定手段と、閾値決定手段で決定された各入力属性の閾値に基づいて、「入力属性が閾値以下であれば第2データ群に含まれるデータである」という相関ルールの確からしさを表す第1のルール評価値と、「入力属性が閾値を超えていれば第2データ群に含まれるデータである」という相関ルールの確からしさを表す第2のルール評価値とを各入力属性について演算する第2の評価手段と、全ての入力属性に関する相関ルールのうちで最も高いルール評価値を持つ相関ルールの入力属性条件を示すデータを、第2データ群に対応する出力属性条件の要因を示す情報として抽出する要因抽出手段とを含むようにしてもよい。

20

また、本発明に係るデータ分析方法は、上記の課題を解決するために、前記のデータ分析装置を用いて、複数の入力属性と、出力属性とで構成されるデータの集合である基本データ群を分析対象とし、入力属性と出力属性との因果関係を分析し、因果関係を示す情報を抽出するデータ分析方法であって、上記分類手段により、基本データ群を出力属性に依って第1データ群と第2データ群とに分類する分類ステップと、上記第1の評価手段により、各入力属性の全ての数値について、入力属性がその数値以下であるデータが第1データ群および第2データ群のうち一方に偏っている度合いを表す閾値評価指標を演算する第1の評価ステップと、上記閾値決定手段により、第1の評価ステップで演算された閾値評価指標に基づいて、各入力属性について最大の閾値評価指標を持つ数値を各入力属性の閾値として決定する閾値決定ステップと、上記第2の評価手段により、閾値決定ステップで決定された各入力属性の閾値に基づいて、「入力属性が閾値以下であれば第2データ群に含まれるデータである」という相関ルールの確からしさを表す第1のルール評価値と、「入力属性が閾値を超えていれば第2データ群に含まれるデータである」という相関ルールの確からしさを表す第2のルール評価値とを各入力属性について演算する第2の評価ステップと、上記要因抽出手段により、全ての入力属性に関する相関ルールのうちで最も高いルール評価値を持つ相関ルールの入力属性条件を示すデータを、第2データ群に対応する出力属性条件の要因を示す情報として抽出する要因抽出ステップとを含むようにしてもよい。

30

40

また、本発明に係るデータ分析プログラムは、上記の課題を解決するために、コンピュータを、基本データ群を出力属性に依って第1データ群と第2データ群とに分類する分類手段、各入力属性の全ての数値について、入力属性がその数値以下であるデータが第1デ

50

ータ群および第2データ群のうちの一方向に偏っている度合いを表す閾値評価指標を演算する第1の評価手段、第1の評価手段で演算された閾値評価指標に基づいて、各入力属性について最大の閾値評価指標を持つ数値を各

入力属性の閾値として決定する閾値決定手段、閾値決定手段で決定された各入力属性の閾値に基づいて、「入力属性が閾値以下であれば第2データ群に含まれるデータである」という相関ルールの確からしさを表す第1のルール評価値と、「入力属性が閾値を超えていれば第2データ群に含まれるデータである」という相関ルールの確からしさを表す第2のルール評価値とを各入力属性について演算する第2の評価手段、および全ての入力属性に関する相関ルールのうちで最も高いルール評価値を持つ相関ルールの入力属性条件を示すデータを、第2データ群に対応する出力属性条件の要因を示す情報として抽出する要因抽出手段として機能させるためのデータ分析プログラムであってもよい。

10

また、本発明に係るデータ分析装置は、上記要因抽出手段で抽出された入力属性条件に基づいて、基本データ群を、上記入力属性条件を満たす要因データ群と上記入力属性条件を満たさない他データ群とに分割し、分類されたデータ群のうち少なくとも一方を新たな基本データ群として分類手段に送る分割手段をさらに含み、分類手段による処理、第1の評価手段による処理、閾値決定手段による処理、第2の評価手段による処理、要因抽出手段による処理、および分割手段による処理からなる一連の処理が繰り返し実行されるようになっていてもよい。

【図面の簡単な説明】

【0112】

20

【図1】本発明の一実施形態に係るデータ分析装置の構成を示すブロック図である。

【図2】本発明の一実施形態に係るデータ分析方法を示すフローチャートである。

【図3】本発明の一実施形態に係るデータ分析装置における頻度累積差演算部7（ステップ5）の出力の一例をグラフで表したもので、入力属性x1と、良品の1-x1頻度累積（A）、不良品の2-x1頻度累積（B）、x1頻度累積差|A-B|との関係を示す。

【図4】本発明の一実施形態に係るデータ分析装置における頻度累積差演算部7（ステップ5）の出力の一例をグラフで表したもので、入力属性x2と、良品の1-x2頻度累積（A）、不良品の2-x2頻度累積（B）、x2頻度累積差|A-B|との関係を示す。

【図5】本発明の一実施形態に係るデータ分析装置における頻度累積差演算部7（ステップ5）の出力の一例をグラフで表したもので、入力属性x3と、良品の1-x3頻度累積（A）、不良品の2-x3頻度累積（B）、x3頻度累積差|A-B|との関係を示す。

30

【図6】本発明の一実施形態に係るデータ分析装置における頻度累積差演算部7（ステップ5）の出力の一例をグラフで表したもので、入力属性x4と、良品の1-x4頻度累積（A）、不良品の2-x4頻度累積（B）、x4頻度累積差|A-B|との関係を示す。

【図7】本発明の一実施形態に係るデータ分析装置における寄与率演算部13（ステップ12）で出力されるデータの一例であり、問題事象である出力属性条件y=2(=Y)に対する入力属性条件「x1>2」および入力属性条件「x2>2」の寄与率を示す。

【図8】本発明の実施形態の入力属性閾値テーブルを、決定木の形式で表現した図である。

【図9】従来の決定木-2を、図7と同じ形式で表現した図である。

40

【図10】従来の決定木-1を表す図である。

【図11】従来の決定木-2のラベル階層構造を表す図であり、(a)はx1属性、(b)はx2属性、(c)はx3属性、(d)はx4属性を示す。

【図12】従来の決定木-2を表す図である。

【符号の説明】

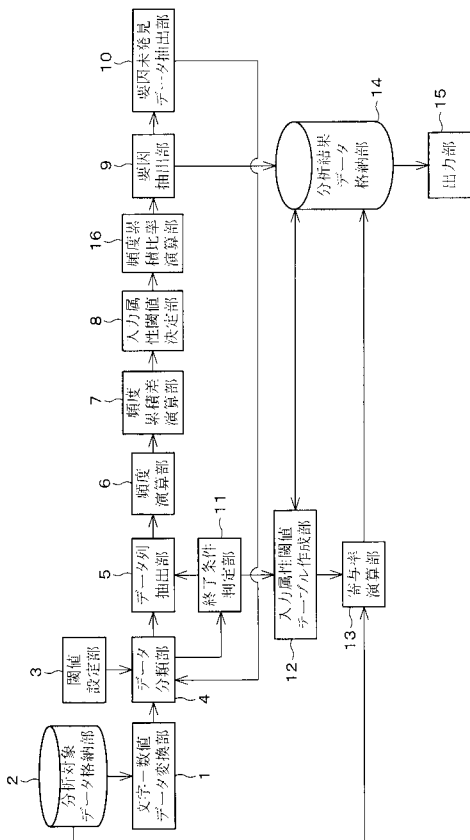
【0113】

- 3 閾値設定部（閾値設定手段）
- 4 データ分類部（分類手段）
- 6 頻度演算部（第1の評価手段、頻度演算手段）
- 7 頻度累積差演算部（第1の評価手段、差分演算手段）

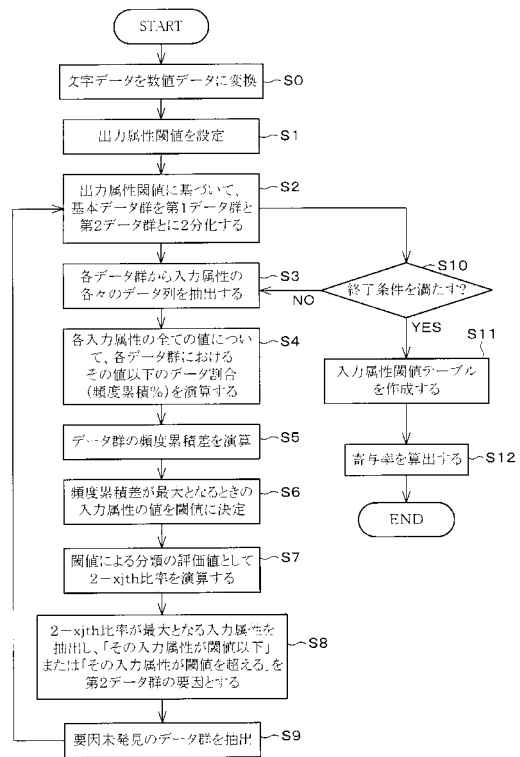
50

- 8 入力属性閾値決定部 (閾値決定手段)
- 9 要因抽出部 (要因抽出手段)
- 10 要因未発見データ抽出部 (分割手段)
- 11 終了条件判定部 (終了条件判定手段)
- 16 頻度累積比率演算部 (第2の評価手段)

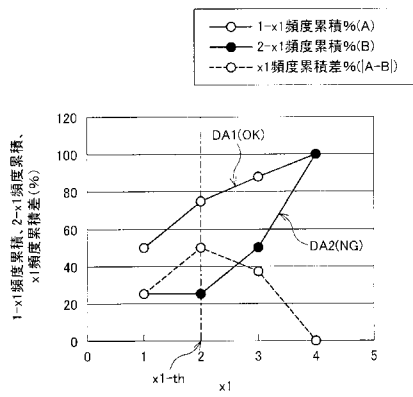
【図1】



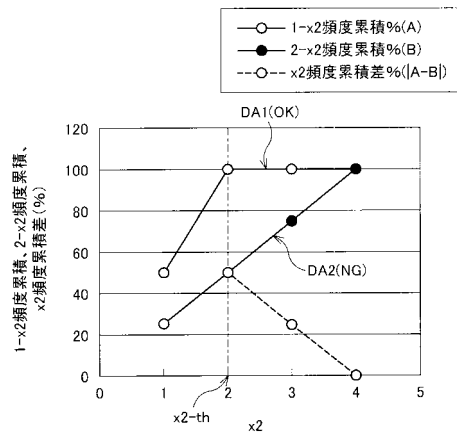
【図2】



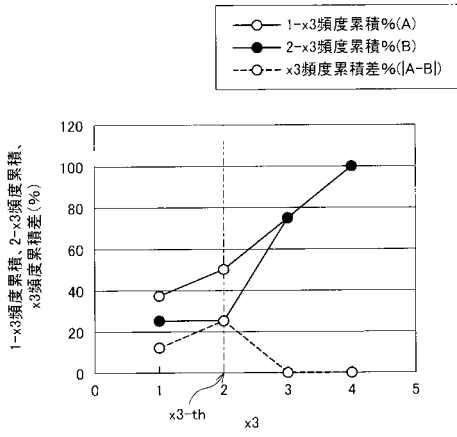
【 図 3 】



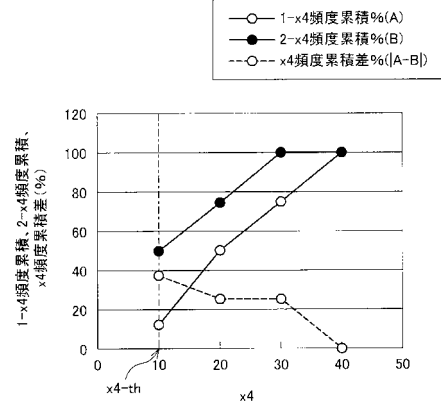
【 図 4 】



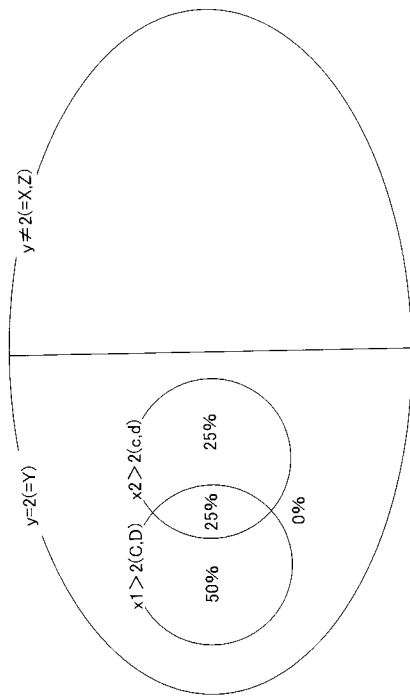
【 図 5 】



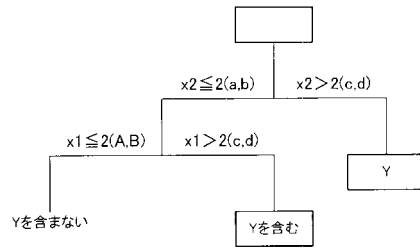
【 図 6 】



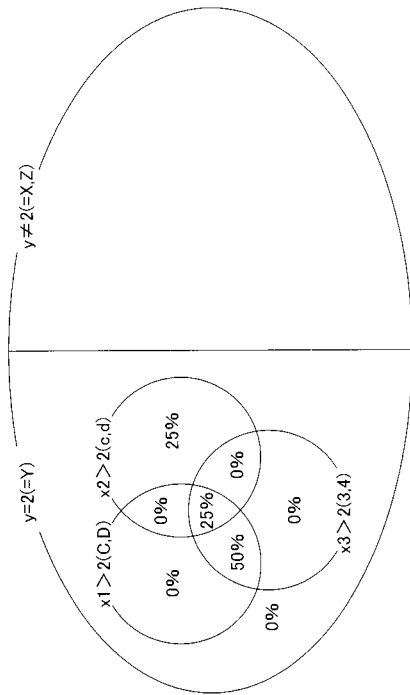
【 図 7 】



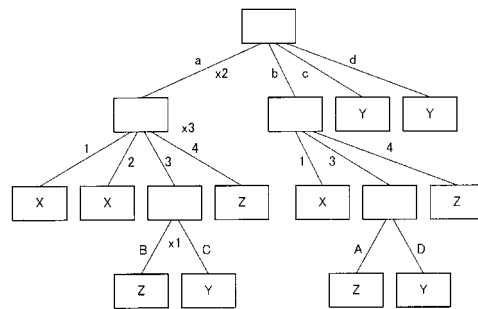
【 図 8 】



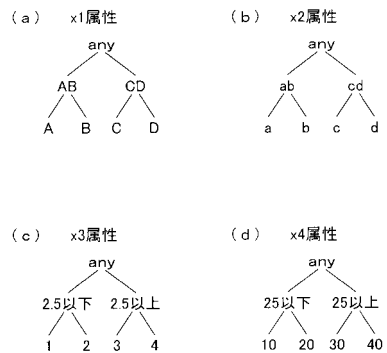
【 図 9 】



【 図 10 】



【 図 11 】



【 図 1 2 】

