

(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2014-13601

(P2014-13601A)

(43) 公開日 平成26年1月23日(2014.1.23)

(51) Int.Cl.

G06F 17/30 (2006.01)

F I

G06F 17/30 210D

テーマコード (参考)

審査請求 有 請求項の数 18 O L (全 27 頁)

(21) 出願番号 特願2013-190084 (P2013-190084)
 (22) 出願日 平成25年9月13日 (2013. 9. 13)
 (62) 分割の表示 特願2010-539868 (P2010-539868)
 の分割
 原出願日 平成20年12月19日 (2008. 12. 19)
 (31) 優先権主張番号 61/015, 973
 (32) 優先日 平成19年12月21日 (2007. 12. 21)
 (33) 優先権主張国 米国 (US)

(71) 出願人 592053963
 エム ケー エス インストルメンツ イ
 ンコーポレーテッド
 MKS INSTRUMENTS, INC
 ORPORATED
 アメリカ合衆国マサチューセッツ州018
 10, アンドーバー, テック・ドライブ
 2, スイート 201
 (74) 代理人 100140109
 弁理士 小野 新次郎
 (74) 代理人 100075270
 弁理士 小林 泰
 (74) 代理人 100101373
 弁理士 竹内 茂雄

最終頁に続く

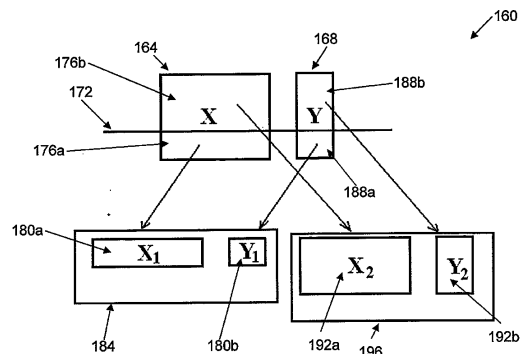
(54) 【発明の名称】 部分的最小二乗分析 (P L S - ツリー) を用いたデータの階層編成

(57) 【要約】 (修正有)

【課題】 比較的短い処理時間で大量のデータを区分 (クラスタリング) する。

【解決手段】 方法は、第1データ・マトリクスおよび第2データ・マトリクスを供給することを伴い、第1および第2データ・マトリクスの各々は、1つ又は複数の変数と、複数のデータ点とを含む。また、本方法は、部分的最小二乗 (P L S) 分析または直交 P L S (O P L S) 分析を用いて第1データ・マトリクスから第1スコアを決定すること、ならびに第1および第2データ・マトリクスを (例えば、行単位で) 第1グループおよび第2グループに区分することを伴い、ソートした第1スコア、第1データ・マトリクスの分散、ならびに第1および第2データ・マトリクスの分散に関連する第1および第2グループの分散に基づいて区分する。

【選択図】 図 1 B



【特許請求の範囲】

【請求項 1】

コンピュータにより実行される方法であって、

第 1 データ・マトリクスおよび第 2 データ・マトリクスを供給する、コンピュータにより実行されるステップであって、前記第 1 および第 2 データ・マトリクスの各々が、1 つ又は複数の変数（マトリクス列）と、複数のデータ点（マトリクス行）とを含む、供給するステップと、

部分的最小二乗（PLS）分析または直交 PLS（OPLS）分析を用いて、前記第 1 データ・マトリクスから第 1 スコアを決定する、コンピュータにより実行される、決定するステップと、

10

第 1 スコアの区分値に基づいて、前記第 1 データ・マトリクスおよび第 2 データ・マトリクスを第 1 グループおよび第 2 グループに行単位に区分する、コンピュータにより実行されるステップであって、該第 1 スコアの区分値は、 (i) 前記第 1 および第 2 グループにおける前記第 1 データ・マトリクスの前記第 1 スコアの分散、 (ii) 前記第 1 および第 2 グループにおける前記第 2 データ・マトリクスの分散、及び (iii) 区分後の前記第 1 グループのサイズと第 2 のグループのサイズとの差に関連する損失関数の合計を備える関数を最小化することによって求められる、区分するステップと、

を備えている、方法。

【請求項 2】

請求項 1 記載の方法であって、区分するステップは、前記 PLS 分析または OPLS 分析を用いることにより生成される前記第 1 スコアの分散と、前記第 2 データ・マトリクスの分散との間の関係を表すパラメータ u を最小化するステップを含む、方法。

20

【請求項 3】

請求項 1 記載の方法であって、前記第 1 データ・マトリクスは、プロセス・データを表すデータを収容する、方法。

【請求項 4】

請求項 1 記載の方法であって、前記第 2 データ・マトリクスは、収量データ、品質データ、またはその組み合わせを表すデータを収容する、方法。

【請求項 5】

請求項 1 記載の方法であって、前記第 1 データ・マトリクスは、対象の分子または高分子の構造的変動に関連のある測定データまたは計算データを表すデータを収容する、方法。

30

【請求項 6】

請求項 1 記載の方法であって、前記第 2 データ・マトリクスは、同じ分子または高分子の生物学的データを表すデータを収容する、方法。

【請求項 7】

請求項 1 記載の方法であって、前記第 1 グループは、前記第 1 および第 2 データ・マトリクスを前記第 1 および第 2 グループに行単位に区分した結果各々得られた、第 3 データ・マトリクスおよび第 4 データ・マトリクスを含み、前記方法は、更に、

第 2 部分的最小二乗（PLS）分析または OPLS 分析を用いて、前記第 3 データ・マトリクスから第 2 スコアを決定する、コンピュータにより実行されるステップと、

40

前記第 2 スコアの第 2 区分値に基づいて、前記第 3 データ・マトリクスおよび第 4 データ・マトリクスを行単位に前記第 3 グループおよび第 4 グループに区分する、コンピュータにより実行されるステップであって、該第 2 区分値は、前記第 3 データ・マトリクスおよび前記第 4 データ・マトリクスの前記第 2 スコアに関連する前記第 3 グループ、および第 4 グループ内の分散を備える第 2 関数を最適化することによって求められる、区分するステップと、

を備えている、方法。

【請求項 8】

請求項 7 記載の方法であって、前記第 2 グループは、第 5 データ・マトリクスおよび第

50

6 データ・マトリクスを含み、前記方法は、更に、

前記第 2 グループが閾値数よりも多いデータ点を含む場合、第 3 部分的最小二乗 (P L S) 分析または O P L S 分析を用いて、前記第 5 マトリクスから第 3 スコアを決定する、コンピュータにより実行される、決定するステップと、

前記第 3 スコアの第 3 区分値に基づいて、前記第 5 データ・マトリクスおよび第 6 データ・マトリクスを行単位に前記第 5 グループおよび第 6 グループに区分する、コンピュータにより実行されるステップであって、該第 3 区分値は、前記第 5 データ・マトリクスおよび前記第 6 データ・マトリクスの前記第 3 スコアに関連する前記第 5 グループおよび第 6 グループ内の分散を備える第 3 関数を最適化することによって求められる、区分するステップと、を備えている、方法。

10

【請求項 9】

請求項 8 記載の方法であって、更に、前記第 1、第 2、第 3、第 4、第 5、または第 6 グループを階層的に表示するステップを備えている、方法。

【請求項 10】

請求項 8 記載の方法であって、更に、前記第 2 グループが閾値数未満のデータ点を含む場合、前記第 2 グループの区分を終了するステップを備えている、方法。

【請求項 11】

請求項 8 記載の方法であって、更に、前記第 2 スコアおよび前記第 2 データ・マトリクスの組み合わせた分散が、前記第 2 グループを前記第 5 および第 6 グループに区分したときに減少しない場合、前記第 2 グループの区分を終了するステップを備えている、方法。

20

【請求項 12】

請求項 8 記載の方法であって、更に、前記第 1 および第 2 データ・マトリクスと関連のある以前の区分値が所定の閾値に等しいかまたはこれを超える場合、前記第 2 グループの区分を終了するステップを備えている、方法。

【請求項 13】

請求項 12 記載の方法であって、前記所定の閾値は、デンドログラムにおける階層レベルの最大数を表す制限値である、方法。

【請求項 14】

請求項 1 記載の方法であって、更に、前記第 1 データ・マトリクスおよび前記第 2 データ・マトリクスを表示するグラフ上において、前記第 1 グループおよび前記第 2 グループを表示するステップを備えている、方法。

30

【請求項 15】

請求項 1 記載の方法であって、四分位 (inter-quartiles) を用いて前記分散を計算する、方法。

【請求項 16】

請求項 1 記載の方法であって、更に、前記第 1 データ・マトリクスおよび前記第 2 データ・マトリクスを、第 1 セットのデータ行を備える第 1 グループおよび第 2 セットのデータ行を備える第 2 グループに行単位に区分するステップを備えている、方法。

【請求項 17】

データを階層的に編成するシステムであって、

40

(a) メモリであって、

(a 1) 第 1 データ・マトリクスおよび第 2 データ・マトリクスを含むデータ構造を含む、メモリと、

(b) 前記メモリに動作的に結合されているプロセッサであって、

(b 1) 前記第 1 データ・マトリクスの部分的最小二乗分析または O P L S 分析に基づいて第 1 スコアを決定するモジュールと、

(b 2) 前記第 1 スコアの区分値に基づいて、第 1 グループおよび第 2 グループを発生するために前記第 1 および第 2 データ・マトリクスを区分するモジュールであって、該第 1 スコアの区分値は、(i) 前記第 1 グループおよび第 2 グループにおける前記第 1 データ・マトリクスの前記第 1 スコアの分散、(i i) 前記第 1 グループおよび第 2 グループ

50

ブにおける前記第2のデータ・マトリクスの分散、及び(i i i)区分後の前記第1グループのサイズと第2グループのサイズとの差に関連する損失関数の合計を備える関数を最小化することによって求められる、モジュールと、

を備えている、プロセッサと、
(c)前記第1および第2グループ、ならびに前記第1および第2グループの前記第1および第2データ・マトリクスに対する関連を表示するために、前記プロセッサに動作的に結合されているディスプレイと、
を備えている、システム。

【請求項18】

データを分析するシステムであって、

メモリから第1データ・マトリクスおよび第2データ・マトリクスを読み出すデータ読み出し手段であって、前記第1および第2データ・マトリクスの各々が1つ又は複数のデータ点を含む、データ読み出し手段と、

部分的最小二乗(PLS)分析またはOPLS分析を用いて、前記第1データ・マトリクスから第1スコアを決定するデータ分析手段と、

第1スコアの区分値に基づいて、前記第1データ・マトリクスおよび第2データ・マトリクスを第1グループおよび第2グループに分割するデータ区分手段であって、該区分値は、(i)前記第1および第2グループにおける前記第1データ・マトリクスの前記第1スコアの分散、(ii)前記第1および第2グループにおける前記第2データ・マトリクスの分散、及び(iii)区分後の前記第1グループのサイズと第2のグループのサイズとの差に関連する損失関数の合計を備える関数を最小化することによって求められる、データ区分手段と、
を備えている、システム。

【発明の詳細な説明】

【技術分野】

【0001】

[0001] 本発明は、一般的には、大きなデータ集合の分析に関し、更に特定すれば、端的にはPLS-ツリーと呼ばれている、部分的最小二乗分析を用いたデータの階層的編成および分析に関する。

【背景技術】

【0002】

[0002] 多くの産業において、製造ならびに研究および開発の双方において、非常に大きなデータ集合が収集される。

【0003】

[0003] 半導体デバイス製造業界では、デバイス製造業者は、一層優れたおよび/または高速のプロセスおよびハードウェア構成を設計するためにプロセス・ツール製造業者を拠り所にすることによって、許容範囲を一層狭めたプロセスおよび材料仕様に移行するように管理している。しかしながら、デバイスの外形がナノメートル等級(nanometer scale)まで縮小しているため、製造プロセスが増々複雑になり、プロセスおよび材料の仕様を満たすのも一層難しくなっている。

【0004】

[0004] 現在の半導体製造において用いられている典型的なプロセス・ツールは、数千個のプロセス変数の集合によって記述することができる。これらの変数は、一般に、製造プロセスの物理的パラメータ、および/または製造プロセスにおいて用いられるツールに關係付けられている。場合によっては、これら数千個の変数の内、数百個の変数が動的であることもある(例えば、製造プロセスにおいて、または製造プロセス間で時間的に変化する)。動的変数には、例えば、気体流量、気体圧力、配給電力、電流、電圧、および温度は、例えば、特定の処理方法(recipe)、処理ステップの全体的なシーケンスにおける特定のステップまたは一連のステップ、製造プロセスの間に起こる誤りおよび故障、または特定のツールまたはチェンバの使用に基づくパラメータ値の変化(例えば、「ドリフト」と

10

20

30

40

50

呼ばれる)に基づいて変化する。

【0005】

[0005] プロセス変数は、収量(yield)変数または応答変数と関係付けられることが多い。プロセス変数は、予測子と考えることができる。即ち、変数間の基礎となる関係に基づいて収量変数を示すと考えることができる。プロセスおよび収量変数を示すデータは、製造プロセスの間に測定され、リアル・タイムの分析または後の分析のために格納される。

【0006】

[0006] 同様に、薬品およびバイオテク生産では、U.S. Food and Drug Administration (米食品医薬品局)のような監督官庁が、指定された品質概要を中心としてばらつきが非常に少ない高品質の製品を維持するために、製造プロセスに対する厳格な仕様の遵守を要求している。これらの仕様は、プロセス変数のオンライン測定、ならびに、例えば、プロセス・ガス・クロマトグラフィ、近赤外線分光分析、および質量分光分析のような追加の多次元センサ技法を必要とする。理想的には、製造プロセスの間に測定されたデータがリアル・タイムの分析に利用することができ、プロセス状態がプロセス仕様にどの位近いかに関する指示または情報を提供することである。

10

【0007】

[0007] 薬品およびバイオ技術の研究および開発では、数万以上のことも多い、多くの異なる分子が、新たな薬品を発見し最適化するプロセスの間に調べられる。多くの異なる物理的および生物的特性が、分子毎(例えば、潜在的な薬品候補)に測定および/または計算され、多くの理論的構造関係特性が分子毎に計算される。分子毎に決定された変数値の総数は数千を越える(例えば、2,000個の変数値よりも多い)こともしばしばである。開発プロセスの一部には、一方では生物的特性と、他方では物理的、化学的、そして理論的に計算された構造関係特性との間における関係を発見することを含む。これらの関係を理解することによって、研究者は将来性のある分子の化学構造を変更し、生物的特性の概要を改善した新たな分子に向かって進むのに役立つ。

20

【0008】

[0008] 大きなデータ集合では、多くの場合、データを互いにグループ化して、クラスタ化データを得る。データに対して有意な分析を行うために、同質のデータ即ちグループ化されないデータ間の比較が好ましい。したがって、グループ化したデータを同質のサブグループにクラスタ化するためのアルゴリズムが開発されている。

30

【0009】

[0009] グループ化したデータを分析する1つの方法は、データに関する線形回帰分析の変形体(例えば、「回帰ツリー」または「分類および回帰ツリー」即ち「CART」と呼ばれることもある)を用いることである。回帰ツリー分析は、個々のX-変数またはX-変数の組み合わせに基づく一連のデータ分割を伴う。データを分割することができる可能性がある方法の数は、観察する変数の数と共に急激に増加する。この理由のために、回帰ツリーが一般に適しているのは、数個の変数だけを有するデータ集合であり、10から20個よりも多い変数を有するデータ集合の場合、回帰ツリー分析は、部分的に計算コストのために、データ集合を細分するのが一般的である。回帰ツリー分析の結果に基づいて、データをツリーまたは分岐編成にグループ化する。これは、デンドログラム(dendrogram)と呼ばれることもある。

40

【0010】

[0010] 階層データ・クラスタリングの一種に、主成分分析(PCA)に基づくものがある。このような技法は、階層レベル毎に、データ集合をPCA分析の第1主成分軸上に投影することを伴う。こうして、投影されたデータは、第1主成分軸に沿って一次的に整列され、第1主成分軸上の中央位置付近においてデータは区分(partition)される。この種の区分またはクラスタリングは、クラスタ・メンバ間の最大距離が所定の(例えば、ユーザが定めた)閾値を超過するまで、再帰的に反復される。CART分析と同様、PCAに基づく分析も、大きなデータ集合には比較的遅い。更に別の欠点は、PCAに基づく分析が通常X-変数のみを考慮し、得られたデータ関係に対するY-変数の影響を無視する

50

ことである。

【 0 0 1 1 】

[0011] 別の技法に、Y - 変数を2つのランダム・グループに分割する、ランダム、二進(0または1) Y - ベクトル値を伴うものがある。部分的最小二乗(PLS)アルゴリズムを用いて、1 - 成分モデルを用いて新たなY - 変数を予測し、予測したY - 変数がランダムY - 変数値に取って代わる。分析が収束した後、予測したY - 変数を最も近い整数(例えば、0または1のいずれか)に丸め、この丸めたY - 変数を用いて、データをグループに区分する。PCAに基づく分析およびCART分析と同様、この手法は、内部計算用PLSを用いるにも拘わらず、X - 変数のみについて動作する傾向がある。この技法の一例では、二進区分(0または1)の代わりに多数(例えば、3、4、またはそれ以上)の区分についてフレームワークを確立することによって、2つよりも多いクラスタを可能にする。

10

【 0 0 1 2 】

[0012] ニューラル・ネットワーク型分析は、データ分析の別の手法である。しかしながら、ニューラル・ネットワーク型分析は、多くの用途に適するだけ十分に計算が高速にはなっておらず、しかも変数の数が10から20を超えると困難が生ずる。

【 発明の概要 】

【 発明が解決しようとする課題 】

【 0 0 1 3 】

[0013] 以前からの手法の欠点には、多数の変数および変数の組み合わせにおいて多くの潜在的な分割を調査する際の集約的計算およびコストが含まれる。回帰ツリーおよびニューラル・ネットワーク型分析は、変数の数が普通または多い(例えば、約20を超える)場合、困難に直面し困難を生ずる。

20

【 課題を解決するための手段 】

【 0 0 1 4 】

[0014] 本明細書に記載する概念は、部分的最小二乗(「PLS」)手法を用いたデータ分析およびデータのクラスタリングまたはグループ化を伴う。部分的最小二乗手法を用いてデータを分析することによって、反復プロセスにおいて比較的大きなデータ集合を小さな部分集合(グループまたはクラスタとも呼ぶ)に区分する。データを区分する毎に得られるグループは、内部同質性(internal homogeneity)(例えば、クラスタ内における変動が少ない)および最大外部異質性(external heterogeneity)(例えば、他のクラスタ(群)に対して相対的に多い変動)のレベルが高くなる。データ分析および編成に対する部分的最小二乗手法は、大きなデータ集合を同様の観察またはデータ点(例えば、プロセス変数および収量変数を関係付けるデータ点)のクラスタまたはグループに分離でき、以前の手法に伴う計算の集中やコストが不要であるという利点がある。また、部分的最小二乗手法は、分析を補助するサブグループ(クラスタ)におけるプロセス変数と収量変数との間の関係を保存する。

30

【 0 0 1 5 】

[0015] 部分的最小二乗手法は、10,000個よりも多い変数を有するデータ集合を含む、比較的大量の変数を有するデータ集合に対処することができる。更に、部分的最小二乗手法は、データ集合が、多数の共線変数またはプロセス変数と収量変数との間に多数の関係を含む場合でも動作することができ、および/または、例えば、検出器の異常またはデータ格納の問題によってデータ集合からデータが部分的に失われた場合でも動作することができる。部分的最小二乗手法の別の利点は、コンピュータ処理時間が比較的速く、比較的速い計算および/またはグラフ即ちプロット上に階層的に編成したデータの提示が容易であることである。

40

【 0 0 1 6 】

[0016] 部分的最小二乗手法を用いることの別の利点は、Y - 変数(例えば、応答)がクラスタリングおよびクラスタリングする決定に影響を及ぼすことである。例えば、Y - 変数は、明示的に、「分割判断基準」の一部、またはデータをサブグループに区分すべきか

50

そしてどこで区分すべきか判断するパラメータとして用いることができる。別の利点として、部分的最小二乗手法は、二進および連続 Y - 変数ならびに 1 つまたは多数の Y - 変数でも動作可能である。部分的最小二乗分析は、分析はデータ集合全体から始まり、連続的にデータをより小さいグループに分割していくということから、「トップ - ダウン」手法である。トップ - ダウン手法は、1 つのデータ点を有するグループから始まって、グループを組み合わせ (2 データ点のグループにする)、全ての観察点が組み合わされて 1 つのデータ・クラスタになるまで続けられるボトム - アップ手法と対照をなす。ボトム - アップ手法は計算コストがかかる傾向がある。何故なら、データをクラスタリングするか否かの判断が、観察点間距離 (例えば、X - 変数間の類似度)、ならびに観察点 - クラスタおよびクラスタ - クラスタ距離 (例えば、同質性および異質性における相違) を監視しなければならないからである。本明細書に記載する手法の別の利点は、本方法が、データ集合または観察点における逸失またはノイズの多いデータにも拘わらず、有用な結果を生成することである。本発明の実施態様には、四分位時間を用いた分散の計算を特徴とするものもある。

10

20

30

40

50

【0017】

[0017] 部分的最小二乗分析に基づくグループ化、区分、またはクラスタリングの利点は、X - 変数自体の値ではなく、PLS 回帰モデルの X - スコアを、分割またはグループ化判断基準の一部として用いることである。PLS 手法を 1 対のマトリクス、X - 変数の X - マトリクスおよび Y - 変数の Y - マトリクスに適用すると、その結果、一連のデータ分割、グループ化、または区分が行われる。データは、行単位 (例えば、観察点単位) で区分され、PLS モデルによって表されるツリー構造またはデンドログラムとなる。このデンドログラムにおける各ノードは、特定のグループまたはクラスタにおけるデータの PLS モデルを表す。

【0018】

[0018] 一般に、1 つのデータ集合またはクラスタを 2 つ (以上) に分割するには、何らかの区分値または位置を決定する。例えば、第 1 マトリクス (X - マトリクス) における X - 変数について第 1 スコア t_1 を計算し、クラスタの観察点をこのスコア t_1 に沿ってソートする。次いで、(a) X - マトリクスの分散、(b) Y - マトリクスの分散、および (c) 後続の各データ集合における観察点の数と関連のある関数 (例えば、損失関数) を含む、数個のファクタ (factor) の加重組み合わせの改良に基づいて区分の位置を決定する。この関数は、データ集合を、実質的に等しくないまたは均衡が取れていないデータ量を有する 2 つのサブグループに分割し難くするファクタと考えることができる。実施形態によっては、サブグループのクロス確認を用いて、デンドログラムの分岐を終了することもある (例えば、更なるサブグループ化が不要であると判断するために)。実施形態によっては、ユーザが PLS - ツリーにおける最大レイヤ数を指定することもあり、典型的な値は 4 または 5 である。

【0019】

[0019] 部分的最小二乗手法は、種々のデータ集合に適用することができる。例えば、本明細書において記載する概念の検査は、プロセス・データ、定量的構造 - 活動関係 (QSAR) データ集合、およびハイパー・スペクトラル画像データ (hyper-spectral image data) に対して実行されたことがある。

【0020】

[0020] 概して言えば、一態様において、比較的短い処理時間で大量のデータを区分するコンピュータ実装方法およびシステムを提供する。本方法は、第 1 データ・マトリクスおよび第 2 データ・マトリクスを供給することを伴う。第 1 および第 2 データ・マトリクスの各々は、1 つ又は複数の変数 (例えば、マトリクス列) と、複数のデータ点 (例えば、マトリクス行) とを含む。また、本方法は、部分的最小二乗 (PLS) 分析または直交 PLS (OPLS) 分析を用いて第 1 データ・マトリクスから第 1 スコアを決定すること、ならびに第 1 および第 2 データ・マトリクスを第 1 データ・グループおよび第 2 データ・グループに区分することを伴い、第 1 データ・マトリクスの第 1 スコア、第 1 データ・マ

トリクスの分散、ならびに第1および第2データ・マトリクスの分散に関連する第1および第2グループの分散に基づいて区分する。一実施形態では、第1および第2データ・マトリクスの各々は、1つ異常のマトリクス列および複数のマトリクス行を含む。

【0021】

[0021] 実施形態によっては、区分するステップは、第1および第2データ・マトリクスを行単位に区分することを伴う。また、区分するステップは、第1 PLSまたはOPLSスコアの分散と、第2データ・マトリクスの変動との間の関係を表すパラメータを最小化するステップを伴うことができる。区分するステップは、第1および第2データ・グループ間における統計的差異を最大化するステップを伴うことができ、第1データ・マトリクスの第1 PLSまたはOPLSスコアの分散、各グループの第2データ・マトリクスの分散、ならびに区分後において第1および第2グループに残っているデータ・マトリクスのサイズ（例えば、サイズの均衡）に関する関数に基づいて、統計的差異を計算する。実施形態によっては、区分するステップは、第1データ・マトリクスの第1スコアの変動、第1および第2データ・グループの各々における第2データ・マトリクスの分散、ならびに区分後に第1および第2データ・グループに残っているデータ（例えば、サイズの均衡）に関する関数を最小化することを伴うこともある。実施形態によっては、PLL-ツリーにおけるレイヤ数がユーザ指定の最大値に達したときに、区分が終了する。

10

【0022】

[0022] 第1データ・マトリクスは、例えば、半導体または薬品および/またはバイオ技術製造プロセスからのプロセス・データを表すデータを収容することができる。更に、第1データ・マトリクスは、例えば、薬品またはバイオ技術研究開発における薬品開発プロジェクトにおいて研究される分子または高分子というような、対象の分子または高分子の構造的変動に関連のあるあるいはこれらを記述する測定データまたは計算データを表すデータを収容することができる。第2データ・マトリクスは、プロセス収量データ、プロセス品質データ、またはその組み合わせを表すデータを収容することができる。別の例では、第2データ・マトリクスは、同じ分子または高分子の生物学的データを表すデータを収容することができる。

20

【0023】

[0023] 実施形態によっては、第1データ・グループは、第1および第2データ・マトリクスを第1および第2グループに行単位に区分した結果各々得られた、第3データ・マトリクスおよび第4データ・マトリクスを含むことがある。このような実施形態では、第3および第4データ・マトリクスの第2部分的最小二乗(PLS)分析またはOPLS分析を用いて、第3データ・マトリクスから第2スコアを決定するステップと、第3データ・マトリクスの第2スコア、第3データ・マトリクスの分散、ならびに第3および第4データ・マトリクスの分散に関連する第3および第4グループにおける分散に基づいて、第3および第4データ・マトリクスを区分する（例えば、行単位の区分）ステップとを伴う。第2データ・グループは、第5データ・マトリクスおよび第6データ・マトリクスを含み、このような実施形態では、本方法は、更に、第2データ・グループが閾値数よりも多いデータ点を含む場合、第3部分的最小二乗(PLS)分析またはOPLS分析を用いて、第5マトリクスから第3スコアを決定するステップと、第5データ・マトリクスの第3スコア、第5データ・マトリクスの分散、ならびに第5および第6データ・マトリクスにおける分散に関連する第5および第6グループにおける分散に基づいて、第5および第6データ・マトリクスを行単位に区分するステップとを伴う。

30

40

【0024】

[0024] 一部の実施形態は、第1、第2、第3、第4、第5、または第6グループを階層的に表示するステップを伴う。更に、本方法は、更に、グループが閾値数未満のデータ点を含む場合、データ・グループの区分を終了するステップを伴うことができる。また、本方法は、更に、第2スコアおよび第2データ・マトリクスの組み合わせた分散が、データ・グループをサブグループに区分したときに減少しない場合、データ・グループの区分を終了するステップを伴うこともできる。実施形態によっては、本方法は、第1および第2

50

データ・マトリクスと関連のある以前の区分が所定の閾値に等しいかまたはこれを超える場合、第2グループの区分を終了するステップを伴う場合もある。所定の閾値は、デンドログラムにおける階層レベルの最大数を表す制限値である。

【0025】

[0025] 実施形態によっては、本方法は、第1データ・マトリクスおよび第2データ・マトリクスを表示するグラフ上において、第1データ・グループまたは第2データ・グループを識別するステップを伴う場合もある。実施形態によっては、ユーザがPLS-ツリーにおけるレイヤの最大数を指定することもあり、典型的な値は4または5である。

【0026】

[0026] 概して、別の態様では、情報担体に有形的に具現化されているコンピュータ・プログラム・プロダクトであって、コンピュータ・プログラム・プロダクトが、データ処理装置に、多数のステップを実行させるように動作可能な命令を含む。例えば、これらのステップは、第1データ・マトリクスおよび第2データ・マトリクスを受け取るステップであって、第1および第2データ・マトリクスの各々が、1つ又は複数のデータ点を含む、ステップと、第1および第2データ・マトリクスの部分的最小二乗(PLS)分析またはOPLS分析を用いて、第1データ・マトリクスから第1スコアを決定するステップと、第1データ・マトリクスの第1スコア、第1データ・マトリクスの分散、ならびに第1および第2データ・マトリクスの分散に関連する第1および第2データ・グループにおける分散に基づいて、第1および第2データ・マトリクスを行単位に区分するステップとを含むことができる。これらの分散は、実施態様によっては、最初の数成分のそれぞれのPLS-スコアの分散によって表される場合もある。

【0027】

[0027] 更に別の態様では、データを階層的に編成するシステムがある。このシステムはメモリを含む。このメモリは、第1データ・マトリクスおよび第2データ・マトリクスを有するデータ構造を含む。また、本システムはメモリに動作的に結合されているプロセッサも含む。このプロセッサは、部分的に第1データ・マトリクスの部分的最小二乗分析またはOPLS分析に基づいて第1スコアを決定するモジュールと、第1グループおよび第2グループを発生するために第1および第2データ・マトリクスを(例えば、行単位で)区分するモジュールを含む。区分は、第1データ・マトリクスの第1スコア、第1データ・マトリクスの分散、ならびに第1および第2データ・マトリクスの分散に関連する第1および第2グループにおける分散に部分的に基づく。また、本システムは、第1および第2データ・グループ、ならびに第1および第2データ・グループの第1および第2データ・マトリクスに対する関連を表示するために、プロセッサに動作的に結合されているディスプレイも含む。

【0028】

[0028] 別の態様では、データを分析するシステムである。このシステムは、メモリから第1データ・マトリクスおよび第2データ・マトリクス(例えば、データ構造)を読み出すデータ読み出し手段を含む。第1および第2データ・マトリクスの各々は、1つ又は複数のデータ点を含む。本システムは、部分的最小二乗(PLS)分析またはOPLS分析を用いて、第1データ・マトリクスから第1スコアを決定するデータ分析手段を含む。また、本システムは、第1および第2データ・マトリクスを第1データ・グループおよび第2データ・グループに分割するデータ区分手段であって、第1データ・マトリクスの第1スコア、第1マトリクスの分散、ならびに第1および第2データ・マトリクスの分散に関連する第1および第2グループにおける分散に基づいて区分する、データ区分手段を含む。

【0029】

[0029] 実施態様の中には、前述の態様のいずれかを含み、以上の実施形態またはその効果の特徴とするものもある。

【0030】

[0030] これらおよびその他の特徴は、以下の説明および図面を参照することによって

10

20

30

40

50

、一層深く理解されよう。図面は、例示であって、必ずしも同じ拡大率で描かれている訳ではない。本明細書では、製造プロセス、特に、半導体、薬品、またはバイオ技術製造プロセスに関してその概念を記載するが、この概念には追加の用途、例えば、データ・マイニング用途、財務データ分析用途、あるいは多数のデータ点または観察を伴うその他の用途もあることは、当業者には明白であろう。

【図面の簡単な説明】

【0031】

[0031] 以上のおよびその他の目的、特徴、および利点は、添付図面に示す実施形態の、以下の更に特定の説明から明白となる。図面において、同様の参照符号は、異なる図面全てを通じて同じ部分を指す。図面は、必ずしも同じ拡大率で描かれているのではなく、むしろ、実施形態の原理を図示する際には強調が加えられている。

10

【図1A】図1Aは、測定データを示すグラフである。

【図1B】図1Bは、データ区分の前および後における、図1Aのグラフ上に表されたデータを示すブロック図である。

【図2】図2は、階層的にデータを編成し表示するデータ処理システムのブロック図である。

【図3】図3は、部分的最小二乗分析を用いてデータを分析する方法を示すフロー・チャートである。

【図4】図4は、部分的最小二乗ツリー分析の後における階層的編成データを示す分類ツリーである。

20

【図5】図5は、データを表示するためのユーザ・インターフェースの一例である。

【図6】図6は、近似検索を用いるためのアルゴリズムの一例を示すフロー・チャートである。

【発明を実施するための形態】

【0032】

[0039] 図1Aは、測定データ105を示すグラフ100である。データ105は、グラフ100上において複数のデータ点110として表されている。データ点110の各々は、製造プロセスまたはその他の何らかの測定または監視プロセス中に収集または測定されたデータを表す。データ点110は、「観察点」(observation)と呼ばれることもある。グラフ100は、第1軸115と、この第1軸115に垂直な第2軸120とを含む。軸115および120は、プロセス変数(観察可能または予測可能変数と呼ばれることもある)または収量変数(結果または予測変数と呼ばれることもある)を表すことができる。実施形態によっては、これらの軸115および120をX-軸と呼ぶ場合もある。また、軸115および120は、Y-軸と呼ぶこともできる。実施形態によっては、第1軸115および第2軸120の単位は、無次元であるか、または目盛りが振られている。実施形態によっては、グラフ100はX-X空間またはY-Y空間においてデータ105を図示し、グラフ100は、1つ又は複数のデータ・マトリクスの平面(または低次元面)への投影を図示することもできる。これらの軸は、データ・マトリクスにおける変数によって定めることができる。

30

【0033】

[0040] 実施形態によっては、データ点110が、プロセス・データおよび対応する収量データ(例えば、プロセス・データを測定したバッチについての収量データ)を表すデータを順序付けた対の一部である場合もある。実施形態によっては、データ点110は、1つ又は複数のデータ・マトリクスにおけるエントリを表すこともある。例えば、プロセス・データは、第1データ・マトリクスにおけるエントリであることもできる。第1データ・マトリクスをX-マトリクスとも呼ぶ。X-マトリクスは、N行(観察点とも呼ぶ)およびK列(変数とも呼ぶ)を含む、 $N \times K$ マトリクスであることができる。収量データは、第2データ・マトリクスにおけるエントリであることができる。第2データ・マトリクスをY-マトリクスとも呼ぶ。Y-マトリクスは、N行およびM列を含む $N \times M$ マトリクスであることができる。

40

50

【 0 0 3 4 】

[0041] グラフ 1 0 0 は、少なくとも 1 0 , 0 0 0 個のデータ点、そして場合によっては、1 0 , 0 0 0 個よりも遥かに多いデータ点を含むことができる。実施形態によっては、第 1 データ・マトリクスおよび / または第 2 データ・マトリクスにおけるデータは、グラフ 1 0 0 上に表示する前に、前処理を受ける。例えば、グラフ 1 0 0 を作成または表示する前に、第 1 データ・マトリクスおよび第 2 データ・マトリクスをストレージから読み出し、前処理アルゴリズム (図示せず) によってマトリクスにおけるデータの変換、中心合わせ、および / または倍率変換を行う。

【 0 0 3 5 】

[0042] 実施形態によっては、前述の前処理を、第 1 または第 2 データ・マトリクスにおけるデータの統計的分析と関連付ける。例えば、ユーザ (例えば、コンピュータまたは人) が、グラフ 1 0 0 を発生する前に、1 組の倍率パラメータを指定して、データに適用することができる。適した倍率パラメータの特定の値を指定するために、倍率調整ファイル (scaling file) を用いることができる。倍率調整は、後続の処理またはモデル発生にデータを用いる前における、データの一種の前処置 (pre-treatment) または前処理と呼ばれることもある。データ・マトリクスにおける観察点および変数の測定値は、非常に異なる数値範囲を有することが多く、このためにデータにおいて大きな統計的分散が生ずる。部分的最小二乗分析は、一般に、最大共分散投影方法と見なされる。その結果、大きな分散がある変数またはデータは、比較的分散が低い変数よりも、グラフ 1 0 0 上で大きく表現される可能性が高い。比較的分散が大きな変数を第 1 軸 1 1 5 (例えば、X - 軸) に沿って散乱プロット (scatter plot) にプロットし、比較的分散が小さな変数を同じ目盛りの第 2 軸 1 2 0 に沿って散乱プロットにプロットすると、散乱が大きな変数における拡散 (spread) が分散が小さな変数における拡散を支配することがあり得る。矯正手段として、双方の変数のデータ (および軸) の倍率を調整することができる。これらの変数を倍率調整することにより、双方の変数が特定のデータ・モデルに寄与することができる。

【 0 0 3 6 】

[0043] 軸 1 1 5 および 1 2 0 の双方に相対的にまたは近似的に等しい重みを与えるために、データ値を標準化、倍率調整、または重み付けする。これによって、X - マトリクスのエントリおよび Y - マトリクスのエントリ (または変数) がモデルにほぼ等しく寄与することを促進する。倍率調整プロセスは、所定の判断基準にしたがって変数空間における座標軸の長さを規制することを伴う (例えば、各座標軸の長さを同じ分散に設定する) 。データを倍率調整する共通の技法に「単位分散」、「UV」倍率調整、または「自動倍率調整」と呼ばれるものがある。単位分散倍率調整は、データ集合からの特定の対して標準偏差 (例えば、) を計算することを伴う。倍率調整重みは、標準偏差の逆数 (例えば、 $w = 1 /$) として計算する。変数の各値に倍率調整重みを乗算して、倍率調整変数を決定する。データ・マトリクスにおける変数の全ての倍率を調整した後、座標軸 1 1 5 および 1 2 0 の各々は単位分散を有する。

【 0 0 3 7 】

[0044] 実施形態によっては、ユーザが特定の対して変数 (例えば、ノイズが多い変数または関連のない変数) の値を減じたり、またはある種の変数のグラフ 1 0 0 に対する寄与を増大させたいこともあり得る。ユーザは、特定のデータ集合についてこの目的を達成するために、倍率調整重みを修正することができる (例えば、このために分散も修正することができる) 。また、変数の変換を用いて、その変数に与える分布を対称に近付けることも多い。例えば、対数変換、負対数倍率調整 (negative logarithm scaling)、ロジット倍率調整 (log-it scaling)、二乗根倍率調整、第 4 根倍率調整、逆倍率調整、またはべき変換倍率調整 (power transformation scaling) を用いることができる。

【 0 0 3 8 】

[0045] 同様に、マトリクス内にあるデータは、倍率調整した座標系の原点 0 以外のある点 (図示せず) を中心にして配することもできる。これを行う場合、必要に応じて、マトリクス要素を他の点を中心にして配するために、中央値をマトリクス列の各々に加算する

または中央値をマトリクス列の各々から減算することができる。中心合わせおよび倍率調整は双方共、グラフ100を発生する際の計算上の要求および/またはデータの部分的最小二乗分析の計算上の要求を減少させることができる。または、中心合わせおよび倍率調整は、データの解釈、および結果的に得られるパラメータまたは解釈モデル(interpretive model)を使用し易くする。

【0039】

[0046] 第1および第2マトリクスの中にあるデータがインポート、および/または中心合わせ、変形、あるいは倍率調整された場合、部分的最小二乗アルゴリズムをそのデータに適用して、t1スコアを判定する。実施形態によっては、部分的最小二乗アルゴリズムは、データの直交部分的最小二乗分析に基づき、t1スコアはデータのこのOPLS分析に基づくこともある。t1スコアは、データのクラスタを近似し、第2マトリクスにおけるデータに相関付けられた、X-空間(例えば、グラフ100上)にあるラインに対応する。部分的最小二乗成分に沿った座標は、個々のデータ点即ち観察点についてのt1スコアを定義または決定する。X-空間において蓄積した観察点についてのt1スコアは、t1スコア・ベクトルを定義または決定する。t1スコア・ベクトルは新たな変数と見なすことができる。

10

【0040】

[0047] t1スコアは、グラフ100上におけるライン125を表す(例えば、個々のt1スコアの累積、またはt1スコア・ベクトルを表す)。t1スコア(例えば、ライン125)に対して垂直なライン145を用いて、グラフ100を2つのセクション130および135に分割または区分する。セクション135は、ライン145よりも下にある観察点即ちデータ点を表し、セクション130は、ライン145よりも上にある観察点即ちデータ点を表す。グラフ100上にあるデータ105を、t1スコア・ベクトル(例えば、ライン125)に沿ってソートする。ライン125に沿ったt1スコア毎に計算を行う。

20

【0041】

[0048] t1スコア・ベクトルを決定した後、ライン125上にある点(例えば、ライン145に沿った)の各分割値を式1によって評価する。

$$u = (1-b) \cdot \{a[V(t_{11})+V(t_{12})]/V(t_1)+(1-a)[V(y_1)+V(y_2)]/V(y)\} + b \cdot F(n_1, n_2) \quad \text{式1}$$

ここで、

30

u = 最小化すべきパラメータ、

a = ユーザが調節可能なパラメータであり、通例、0と1との間、

b = ユーザが調節可能なパラメータであり、通例、0と1との間、

V = 特定のマトリクスまたはベクトル内における分散

t_{1i} = ライン125上におけるi番目の座標、例えば、i番目の観察点についてのt1スコア値、

y_i = Yマトリクスにおけるi番目の行、例えば、i番目の観察点のY-ベクトル、

n₁ = サブグループ1(例えば、セクション130内にある)の中にあるデータ点即ち観察点の数、

n₂ = サブグループ2(例えば、セクション131内にある)の中にあるデータ点即ち観察点の数、

40

F = n₁およびn₂を関係付ける関数。データを、n₁およびn₂についてほぼ同様の値を有するサブグループに区分し易くするために用いられる。

【0042】

[0049] 式1は、定性的には、Xマトリクス(例えば、t1スコア)における分散と、Y-マトリクスにおける分散と、各潜在的なサブグループまたはサブパーティションにおけるデータ量との間の関係と考えることができる。「u」の値は、通例、X-スコアt1の分散、Y-マトリクスの分散、および連続する各サブグループにおけるデータ量と関連付けられている関数(F(n₁, n₂))の組み合わせの全体的な改良によって、最小化される(つまり、t1スコアに沿った区分が最適化される)。例えば、関数F(n₁, n₂

50

)は、結果的に得られる各サブグループにおいてほぼ等しい数の観察点(例えば、X - 変数)が得られ易くする損失関数と考えることができる。実施形態によっては、この関数 $F(n_1, n_2)$ は式2によって与えられる。

【0043】

【数1】

$$F = \frac{(n_1 - n_2)^2}{(n_1 + n_2)^2}$$

式2

10

【0044】

[0050] 当業者には、他の損失関数も明白であろう。実施形態によっては、ユーザ調節可能パラメータ a が、スコア t_1 および Y - 変数を関係付ける場合もある。例えば、 a の値が0に近づく程、スコア1に起因するウェイトが大きくなる。 a の値が1に近づく程、 Y - 変数に起因するウェイトが大きくなる。ユーザ調節可能パラメータ b は、サブグループのサイズに関係する。例えば、 b の値が0に近づく程、区分によって得られるサブグループは、区分後にほぼ等しいサイズとなる可能性が低くなる。 b の値が1に近づく程、区分によって得られるサブグループは、区分後にほぼ等しいサイズとなる可能性が高くなる。実施形態によっては、パラメータ a のデフォルト値は0.3であり、パラメータ b のデフォルト値も0.3である。パラメータ a および b について、他のデフォルト値も可能である。

20

【0045】

[0051] 実施形態によっては、パラメータ b の値は0にすることができる。このような実施形態では、損失関数 F は、最小化すべきパラメータ u には影響を及ぼさない。具体的には、損失関数 F は、それぞれのサブグループにおけるデータ量の増大を促すためまたはデータ量に影響を与えるためには用いられない。実施形態によっては、パラメータ a の値は0にすることができる。このような実施形態では、式1は、区分化が X - 変数自体の値の代わりに PLS - スコアに基づくことを除いて、分類および回帰ツリー ($CART$) 分析に類似する。

30

【0046】

[0052] 実施形態によっては、パラメータ n_{min} を、例えば、パラメータ b の値が0に近いまたは比較的小さいときに、各々比較的少量のデータを収容する比較的多数のクラスタまたはグループが生ずる式1に対する解を防止する境界条件またはパラメータとして指定することができる。パラメータ n_{min} は、関数的に、 $n_{min} = \min(n_1, n_2)$ と表すことができる。 n_{min} の値の一例は5である。 n_{min} には他の値も可能であり、ユーザが選択することができる。実施形態によっては、階層レベルの数(暗示的に、サブグループまたはクラスタの数)をユーザが決定または選択することができる。例えば、ユーザは4つまたは5つの階層レベルを選択することができる。ユーザが階層レベルの数を選択していない場合、デフォルトを指定することができる(例えば、4階層レベル)。

40

【0047】

[0053] 式1におけるパラメータ「 u 」の値を最小化するライン125上の座標140が決定されるおよび/または突き止められる。パラメータ「 u 」を最小化することによって、 t_1 スコア(または t_1 スコア・ベクトル)および第2マトリクス (Y - マトリクス)における変動にしたがって、データ105が区分される。交差座標140においてライン125に垂直なライン145が決定され、グラフ100上に図示される。ライン145はこのグラフをセクション130および135に分割する。セクション130は、ライン145よりも上にあるデータ105を含み、セクション135は、ライン145よりも下にあるデータ105を含む。セクション130は、第1データ・グループを含み、セクショ

50

ン 1 3 5 は第 2 データ・グループを含む。式 1 を最小化した結果、(i) 第 1 データ・グループないにおけるスコア t_1 の分散および第 1 グループの第 2 マトリクスの分散、ならびに (i i) 第 2 データ・グループ内におけるスコア t_1 の分散および第 2 グループの第 2 マトリクスの分散の組み合わせが、 t_1 の特定の値の選択によって最小化される。この組み合わせを最小化することは、変数 t_1 および Y に関して、第 1 データ・グループと第 2 データ・グループとの間において組み合わせ分散 (combined variance) を最大化することと同等である。

【 0 0 4 8 】

[0054] グラフ 1 0 0 をセクション 1 3 0 および 1 3 5 に分割した後、同様の手順を用いてセクション 1 3 0 および 1 3 5 の各々におけるデータを分析することができる。例えば、セクション 1 3 5 におけるデータ点 1 1 0 は、第 3 データ・マトリクス X_1 (例えば、セクション 1 3 5 におけるデータ 1 0 5 の X -マトリクス値を含む)、および第 4 データ・マトリクス Y_1 (例えば、セクション 1 3 5 におけるデータ 1 0 5 の Y -マトリクス値を含む) と見なすことができる。先に論じたのと同様にして、第 3 データ・マトリクスから第 2 t_1 スコアを決定することができる (しかし、セクション 1 3 5 の中にあるデータのみに基づく)。セクション 1 3 5 は、第 2 t_1 スコアに基づいて第 2 ライン (図示せず) に沿って分割または区分することができる。第 2 t_1 スコアを決定した後、式 1 は、第 2 t_1 スコアおよび第 4 データ・マトリクス (例えば、 Y -マトリクス) における変動に関して最小化し、更にセクション 1 3 5 を第 2 垂直ライン (図示せず) に沿って第 1 および第 2 サブグループ (図示せず) に細分することができる。

10

20

【 0 0 4 9 】

[0055] 次いで、セクション 1 3 0 (グループ 2) の類似した分析が続き、セクション 1 3 0 の中にあるデータをサブグループに区分することができる。

【 0 0 5 0 】

[0056] 以上で説明した手順は、グラフ 1 0 0 上のデータ 1 0 5 全てを分析し、増々小さなクラスタ (サブグループ) の階層構造にグループ化し終わるまで、連続するサブグループ毎に継続することができる。実施形態によっては、この区分プロセスは、サブグループが収容するデータ点が閾値データ点数よりも少なくなったとき、またはデータを更に区分しても t_1 スコア・ベクトルにおける相対的分散または Y -マトリクスにおける変動が小さくならないときに、特定のサブグループに対して終了する。データ点の閾値数は、ユーザが選択可能であり、例えば、5 データ点とすることができる。

30

【 0 0 5 1 】

[0057] 実施形態によっては、クラスタ外形 (geometry) の検査から、クラスタまたはサブグループが、第 1 スコア・ベクトル t_1 に平行でない方向に沿って方位付けられていることが示唆されることがある。サブグループが第 1 スコア・ベクトル t_1 に平行に方位付けられていない場合、スコア・ベクトルの組み合わせ (例えば、2 つ、3 つ、またはそれ以上のスコア・ベクトルの組み合わせ) を用いることができる。スコア・ベクトルを組み合わせるには、実施形態によっては、第 3 パラメータ c を導入することもある。第 3 パラメータ c は、通常 - 1 と + 1 との間の値を有する。第 3 パラメータ c は、第 1 スコア t_1 を第 2 スコア t_2 に関係付ける。パラメータ c を用いた、スコア・ベクトル t_1 と t_2 との間における適した関係の一例は、次の通りである。 $\{c \cdot t_2 + (1 - |c|) \cdot t_1\}$ 。他の関係も当業者には明白であろう。パラメータ c とスコア・ベクトル t_1 および t_2 との間における関係から、第 1 スコア・ベクトル t_1 または第 2 スコア・ベクトル t_2 のみに沿った分析ではなく、スコア・ベクトル t_1 および t_2 によって定義される平面において表されるデータの分析ができる。

40

【 0 0 5 2 】

[0058] 実施態様によっては、ユーザが分析対象変数の数を減少させることができる場合もある (例えば、データ選択または前処理とも呼ばれている)。例えば、ある種の変数は、モデルにおける最良の予測変数と強く相関付けることができ、あるいは Y -変数 (例えば、結果変数) とは相関付けられない。データ選択の一例では、データ (例えば、 X -変

50

数)のパラメータを所定値と比較することを伴う。例えば、Yとの相関が所定の百分率(例えば、75%)よりも低いことを表示する変数を、分析の前に、データ集合から排除する。

【0053】

[0059] 図1Bは、データ区分の前および後における、図1Aのグラフ100上に表されたデータを示すブロック図160である。ブロック図160は、第1データ・マトリクス164(Xで示す)および第2データ・マトリクス168(Yで示す)を含む。データ・マトリクス164、168の各々は、1つ又は複数の列(変数とも呼ぶ)および複数の行(観察点とも呼ぶ)を含むことができる。実施形態によっては、第1データ・マトリクス164がプロセス・データを含み、第2データ・マトリクス168が収量データおよび/または製品品質データを含む場合もある。別の実施形態では、第1データ・マトリクス164が測定および計算した物理化学的データおよび/または構造関係データを含み、第2データ・マトリクス168は、例えば、1組の分子または高分子と関連のある生物データを含む。

10

【0054】

[0060] また、ブロック図160は、第1および第2データ・マトリクス164、168を分割する区分ライン172も含む。第1データ・マトリクス164の部分176aは、第1データ・マトリクス164の区分時に、第1データ・グループ184の部分180a(X_1 で示す)となる。同様に、第2データ・マトリクス168の部分188aは、区分時に、第1データ・グループ184の部分180b(Y_1 で示す)となる。第1データ・マトリクス164の部分176bは、第1データ・マトリクス164の区分時に、第2データ・グループ196の部分192a(X_2 で示す)となる。同様に、第2データ・マトリクス168の第2部分188bは、区分時に、第2データ・グループ196の部分192b(Y_2 で示す)となる。したがって、後続の区分は、第1データ・マトリクス164および第2データ・マトリクス168を表す第1データ・グループ184の第1部分180a(X_1)および部分180b(Y_1)について同様に続けられる。第2データ・グループ196も同様に区分することができる。

20

【0055】

[0061] 実施形態によっては、区分は、第1データ・マトリクス164における観察点のスコア、および第2データ・マトリクス168の分散に応じた行単位の区分である。他の区分技法も用いることができる。実施形態によっては、区分が中止または終了するのは、特定のデータ・グループ(例えば、データ・グループ184)がデータ点の閾値数を超えないとき、例えば、観察点の数が少なすぎる場合である。データ点の閾値数は、ユーザによって設定または決定することができる。実施形態によっては、データ・グループの区分を終了するのは、X-マトリクス(t_1)およびY-マトリクスの第1スコア・ベクトルにおける分散が、区分の結果減少しなくなったときである。これらの分散を比較することができる(例えば、 $[V_1 + V_2] / V$ の比または分数として比較することができる。これは、式1からの変数uと呼ばれることもある)。uが1以上である場合、区分プロセスの結果得られたマトリクスにおける分散は、直前のデータ・グループから得られた分散以上となり、区分を終了する。uが1未満である場合、区分プロセスの結果得られたマトリクスにおける分散は、直前のデータ・グループの分散よりも小さく、uが1以上になるまで区分を継続する。

30

40

【0056】

[0062] 図2は、データを階層的に編成し表示するデータ処理システム200のブロック図である。データ処理システム200は、プロセッサ210に結合されているメモリ205を含む。また、データ処理システム200は、プロセッサ210に結合されているディスプレイ215も含む。更に、データ処理システム200は、図示しないその他のコンポーネントまたはモジュール、例えば、データを測定し、収集し、メモリ205に格納するデータ取込モジュール、または多変量統計分析にしたがって収集したデータに基づいてモデルを発生するモデル発生モジュールも含むことができる。データ処理システム200は

50

、現場分析および/またはリアル・タイム分析のために製造設備に設置することができ、あるいは処理後の分析またはデータ・マイニングの用途のために他の場所に設置することもできる。

【 0 0 5 7 】

[0063] メモリ 2 0 5 は、例えば、処理変数に関するデータ (X - マトリクス・データ) および収量変数に関するデータ (Y - マトリクス・データ) のような、例えば、製造プロセスを表すデータを含む。データは、生データとして、データ・モデルまたはテンプレートとして、あるいは前処理済みのデータ (例えば、倍率調整、中心合わせ、および/または変換後) として格納することができる。

【 0 0 5 8 】

[0064] プロセッサ 2 1 0 は、メモリ 2 0 5 と通信するデータ読み出しモジュール 2 2 0 を含む。データ読み出しモジュール 2 2 0 は、分析のためにメモリ 2 0 5 からデータを読み出す。また、プロセッサ 2 1 0 はデータ分析モジュール 2 2 5 およびデータ区分モジュール 2 3 0 も含む。データ分析モジュール 2 2 5 は、データ読み出しモジュール 2 2 0 およびデータ区分モジュール 2 3 0 と通信する。データ区分モジュール 2 3 0 は、データ読み出しモジュール 2 0 5 と通信する。データ読み出しモジュール 2 0 5 は、ディスプレイ 2 1 5 と通信し、読み出したデータのユーザへの表示をし易くする (例えば、図 1 のグラフ 1 0 0 上にあるデータ点 1 1 0)。データ区分モジュール 2 3 0 およびデータ分析モジュール 2 2 5 もディスプレイ 2 1 5 と通信し、データをユーザに表示し易くする。

【 0 0 5 9 】

[0065] データ読み出しモジュール 2 2 0 がメモリ 2 0 5 から特定のデータ集合を読み出した後、データ分析モジュールは、読み出したデータ (たとえば、第 1 X - マトリクスおよび Y - マトリクス) の部分的最小二乗分析 (P L S) または直交 P L S 分析 (O P L S) を実行して、第 1 t 1 スコアを決定する。先に論じたように、t 1 スコアは、読み出したデータを細分するための基準を形成する X - 空間内のラインを表す。データ分析モジュール 2 2 5 は、ディスプレイ 2 1 5 と通信し、分析したデータのユーザへの表示をし易くする (例えば、図 1 のグラフ 1 0 0 上にあるライン 1 2 5 のように)。

【 0 0 6 0 】

[0066] データ分析モジュール 2 2 5 が t 1 スコアを決定し終わると、データ区分モジュール 2 3 0 が t 1 スコアを X - マトリクスおよび Y - マトリクス (例えば、第 1 および第 2 データ・マトリクス) に関して分析し、区分されたグループ内における分散を最小化し、区分されたグループ間の分散を最大化する値を、t 1 について決定する。データ区分モジュール 2 3 0 は、ディスプレイ 2 1 5 と通信し、区分したグループのユーザへの表示をし易くする (例えば、図 1 のグラフ 1 0 0 上におけるライン 1 4 5 およびセクション 1 3 0、1 3 5 として)。

【 0 0 6 1 】

[0067] 実施形態によっては、データ読み出しモジュール 2 2 0、データ分析モジュール 2 2 5、およびデータ区分モジュール 2 3 0 の内 1 つ又は複数が、同じアプリケーション、プロセス、またはプログラムのサブルーチンあるいはサブアルゴリズムである場合もある。実施形態によっては、データ分析モジュール 2 2 5 およびデータ区分モジュール 2 3 0 が同じサブルーチンまたはアルゴリズムの一部であることもある。

【 0 0 6 2 】

[0068] ディスプレイ 2 1 5 は、ユーザ入力デバイス (図示せず)、例えば、ユーザにパラメータを指定させるまたは命令をプロセッサ 2 1 0 に発行させるキーボードまたはマウスを含むことができる。実施形態によっては、ディスプレイがユーザ・インターフェースを含み、ユーザとプロセッサとの間における通信をやり易くする場合もある。例えば、ユーザは、ユーザ・インターフェースを通じて、前述の式 1 からパラメータ「 a 」および「 b 」の値を指定することができ、またはユーザは命令をプロセッサ 2 1 0 に発行し、データ読み出しモジュール 2 2 0 に、分析のために、メモリ 2 0 5 から指定のデータ集合を読み出すように指令することができる。加えて、ユーザは、データをメモリ 2 0 5 から読み

10

20

30

40

50

出す前または後のいずれでも、前処理モジュール（図示せず）によってデータを前処理するために、データの倍率調整、変換、または中心合わせを指定することができる。ユーザが別のコンピュータ・システムまたはプロセッサ（図示せず）である実施形態では、ユーザ・インターフェースは、メモリ205内にあるデータに関して、そしてプロセッサ210がデータをどのように処理するかに関して、システム200にパラメータを指定することができるマシン・マシン・インターフェースとすることができる。

【0063】

[0069] 本発明の特徴および態様を具現化した市販製品の一例に、スウェーデン、UmeaのUmetrics, Inc.が販売するSIMCA-P+（商標）ソフトウェア製品がある。

【0064】

[0070] 図3は、部分的最小二乗分析を用いてデータを文政する方法を示すフロー・チャート300である。ステップ304は、分析対象データを読み出すことを伴う。読み出されたデータは、マトリクス形状あるいはその他の何らかの形態または構造とすることができる（例えば、図1Aおよび図1Bに関して先に論じた第1および第2データ・マトリクス）。実施形態によっては、データはメモリ（例えば、コンピュータ化されたメモリ）から読み出される。また、ステップ304は、フロー・チャート300に示した方法が完了した後に行われる。例えば、ステップ304は、図1Aおよび図1Bに関して先に論じたように第1および第2グループに分割されたデータの後続処理または一連の繰り返し処理における最初のステップを表すことができる。ステップ304は、フロー・チャート300における方法が行われている間に、アレイまたは一時的メモリに格納されているデータからデータを読み出すことを伴うことができる。

10

20

【0065】

[0071] ステップ304においてデータを読み出した後、任意のステップ308においてデータを前処理することができる。このデータの前処理は、前述のように、読み出されたデータ・マトリクスにおけるデータの変換、中心合わせ、および/または倍率調整を伴う。実施形態によっては、前処理はユーザの指定に回答して行われる（例えば、倍率調整重み、閾値のような前処理パラメータを含む、または特定の繰り返しに回答して）。実施形態によっては、ユーザは、デフォルト値または予め選択してある値を用いて前処理を行うべきことを指定する。データを前処理するか否かは、本方法の初期セットアップの間にユーザによって指定されるデフォルト設定とすることができる。ステップ306は、前処理ステップ308が完了したか否かに関連する判断ステップを示す。

30

【0066】

[0072] データを前処理した後（ステップ308）、ステップ312において部分的最小二乗分析または直交部分的最小二乗分析を選択することができる。部分的最小二乗分析の種類は、ユーザによって指定することができ、またはデフォルト設定とすることもできる。部分的最小二乗分析が選択された場合（ステップ316）、PLS分析を用いて第1 t_1 スコアを決定する（ステップ324）。直交部分的最小二乗（OPLS）分析が選択された場合（ステップ320）、直交部分的最小二乗分析を用いて第1 t_1 スコアを決定する（ステップ324）。第1 t_1 スコアは、第1および第2データ・マトリクスによって定義される空間または座標系（例えば、X-Y空間）におけるラインまたはその他の何らかの曲線によって表すことができる。実施形態によっては、PLSおよびOPLS分析の双方を用いることもある。

40

【0067】

[0073] ステップ324において t_1 スコアを決定した後、ソーティング・プロセスを行う（ステップ328）。その後、分割プロセス（ステップ332）を行い、第1および第2データ・マトリクスを2つの部分（パート1および2）に分割する。パート1および2は、 t_1 スコア値を第1および第2データ・マトリクスにおける分散に関係付けるパラメータを、パート1およびパート2におけるデータ値が最小化するか否か判断するための更に別の処理および分析のための、データの一時的細分化とすることができる。第1および第2データ・マトリクスをパート1および2に分割した後（ステップ332）、この区分

50

を先の式 1 によって評価する (ステップ 336)。ステップ 340 において、式 1 からのパラメータ「a」および「b」を読み出す。実施形態によっては、パラメータ「a」および「b」をユーザが選択することもある。ステップ 344 において、プロセッサは、t1 スコア (ステップ 340) によって表されるラインに沿った値ならびにおよびパラメータ「a」および「b」を用いて、式 1 を分析し、パート 1 および 2 がパラメータ「u」の値を最小化するか否か評価する。パート 1 および 2 が「u」の値を最小化しない場合、プロセッサは、t1 ラインに沿った値を評価し続ける (例えば、スコア・ベクトル)。プロセッサは、パラメータ「u」の値を最小化する 1 つまたは複数の値を特定する (ステップ 348)。式 1 からの「u」の値を最小化する t1 スコアの値を選択し、パート 1 および 2 (ステップ 332) をグループ 1 および 2 と置き換える (ステップ 348)。

10

【0068】

[0074] ステップ 344 から 348 においてパラメータ「u」を最小化した結果、グループ 1 および 2 は、各々、式 1 に基づくパラメータ「a」および「b」の特定の値について、t1 スコアの最小グループ内分散 (intra-Group variance) ならびに X - マトリクス・データおよび Y - マトリクス・データ、そして t1 スコアの最大グループ間分散 (inter-Group variance) ならびに X - マトリクス・データおよび Y - マトリクス・データを有する。グループ 1 および 2 は、ステップ 304 において読み出した第 1 および第 2 データ・マトリクスからのデータの部分集合を収容する。グループ 1 は、例えば、図 1 B に示すような、第 3 データ・マトリクス (例えば、グループ 1 に分割された第 1 X - マトリクスからの値を含む X - マトリクス) および第 4 データ・マトリクス (例えば、グループ 1 に分割された第 1 Y - マトリクスからの値を含む Y - マトリクス) を含むというように考えることができる。同様に、グループ 2 は、例えば、図 1 B に示すような、ステップ 304 において読み出した値を含む X - マトリクスおよび Y - マトリクスを含む 2 つのデータ・マトリクスを含むことができる。

20

【0069】

[0075] ステップ 352 において、グループ 1 およびグループ 2 におけるデータ点の数、ならびにパラメータ「u」の値を評価する。グループ 1 (またはグループ 2) の中にあるデータ点の数 n が、データ点の閾値数 $n_{\text{threshold}}$ よりも小さい場合、または「u」の値が所定値 (例えば、1) を超える場合、グループ 1 (またはグループ 2) を終了する (ステップ 360) (例えば、これ以上グループの区分を行わない)。 n がグループ 1 (またはグループ 2) について $n_{\text{threshold}}$ に等しいかまたはこれを超え、更に「u」の値が所定値 (例えば、1) よりも小さい場合、本方法は問い合わせ (query) ステップ 356 に進む。ステップ 356 において、階層レイヤまたはレベルの数を評価し、指定された限度と比較する。階層レベルの数がこの限度未満である場合、本方法はステップ 304 に戻り、グループ 1 におけるデータ・マトリクスをデータ (例えば、ステップ 316、320 において部分的最小二乗または直交部分的最小二乗分析を実行したデータ・マトリクスとして) として用いて、処理を開始する。階層レベルの数が限度に等しい場合、プロセスを終了する (ステップ 360)。その後、グループ 1 と同様にグループ 2 を処理する。フロー・チャート 300 における方法は、PLS - ツリー・レイヤの指定最大数に達するまで、またはグループの区分が終了しており、データ点の閾値数 $n_{\text{threshold}}$ よりも多いデータ点を含むグループがなくなるまで繰り返される。実施形態によっては、ツリー・レイヤの数 $n_{\text{threshold}}$ および / または「u」の値の限度は、ユーザによって指定されることもある。

30

40

【0070】

[0076] 図 4 は、部分的最小二乗 (または OPLS) ツリー分析後における階層的編成データを示す分類ツリー 400 である。分類ツリー 400 は、第 1 データ集合 408 を含む第 1 階層レベル 404 を含む。第 1 階層レベル 404 は、後続の処理が行われる前にメモリ (または最上位の階層レベル) から読み出したデータを示すことができる。実施形態によっては、第 1 階層レベル 404 が、既に処理されているデータを示すこともある (例えば、第 1 データ集合 408 は、区分されているサブグループを表す)。第 1 データ集合 4

50

08は、例えば、図1Bに示すような、第1データ・マトリクス（例えば、X - マトリクス）および第2データ・マトリクス（例えば、Y - マトリクス）を含む。

【0071】

[0077] また、分類ツリー400は、第2階層レベル412も含む。第2階層レベル412は、第1データ・グループ416および第2データ・グループ420を含む。第1データ・グループ416および第2データ・グループ420は、図3のフロー・チャート300に図示した方法にしたがって、そして図1Bに例示したように決定される。例えば、第1グループ416および第2グループ420は、式1においてパラメータ「u」を最小化することによって特定され、パラメータ「u」は、第1集合408の第1データ・マトリクスの第1t1スコア、および第1集合408の第2マトリクスの分散に基づいて決定することができる。第1t1スコアおよびY - マトリクスに関して、第1グループ416は最小の内部変動を保有し、第2グループ420に関してまたは関係して最大化された変動を保有する。第1グループ416は、例えば、図1Bに示すように、第1データ・マトリクスからのX - マトリクス値およびY - マトリクス値を含む。第2グループ420は、例えば、図1Bに示すような区分に基づいて、第1データ・マトリクスからの異なるX - マトリクス値およびY - マトリクス値を含む。

10

【0072】

[0078] 分類ツリー400は、第3階層レベル424も含む。第3階層レベル424は、第1データ・サブグループ428、第2データ・サブグループ432、第3データ・サブグループ436、および第4データ・サブグループ440を含む。第1サブグループ428および第2サブグループ432は、第1データ集合408から第1グループ416を決定したのと同様に決定される（例えば、図1Bおよび図3に関して先に論じたように）。更に特定すれば、第2t1スコアは、部分的最小二乗（またはOLS）分析に基づいて、第1データ・グループ416の中にあるデータから計算される。第2t1スコアは、X - マトリクスおよびY - マトリクスにおける分散と共に、第1グループ416におけるデータに関して式1からのパラメータ「u」を最小化するために用いられる。パラメータ「u」が最小化されたときに、第1サブグループ428および第2サブグループ432が発生する。第3サブグループ436および第4サブグループ440は、第2グループ420の中にあるデータに基づいて、第2グループ420の部分的最小二乗分析に基づいて計算された第3t1スコアを用いて、同様に決定される。

20

30

【0073】

[0079] 分類ツリー400は、第4階層レベル444も含む。第4階層レベル444は、第3階層レベル424の第1サブグループ428から決定される、第1サブサブグループ448および第2サブサブグループ452を含む。第2サブサブグループ452は、第4階層レベル444においては表されていない。これは、第2サブサブグループ452から作ろうとしたサブサブグループにおけるデータ点数が、分析継続に対するデータ点の閾値レベルを超えなかったため（分岐とも呼ぶ）、またはパラメータuについて含まれる値が1.0を超えているからである。また、第4階層レベル444は、更に別のサブサブグループ456a、456b、456c、456dも含む。これらは、第1サブサブグループ448および第2サブサブグループ452と同様に決定される。

40

【0074】

[0080] 尚、階層レベル404、412、424、および444の各々は、直上にある階層レベルからのデータ（例えば、第1データ・マトリクスおよび第2データ・マトリクス）を表すことは言うまでもない。例えば、第2階層レベル412によって表されるデータ・マトリクスは、第1階層レベル404では（例えば、データ集合408では）存在したが、ソートされておらず、分類もされていない形態であった。分類ツリー400における各「分岐」460は、部分的最小二乗（またはOLS）分析を用いて（例えば、図3のフロー・チャートにしたがって）データを分類または区分するプロセスを表す。分岐460は、図1Bからのブロック図160を表す。図示のように、各データ・グループと直前の階層レベルとの間の図4におけるy - 軸に沿った距離は、ほぼ同一である（例えば、

50

y - 軸に沿ったサブサブグループ 4 4 8 とサブグループ 4 2 8 との間の距離は、サブサブグループ 4 5 6 c とサブグループ 4 4 0 との間の距離にほぼ等しい)。実施形態によっては、異なる階層レベル上におけるデータ・グループ間の距離が異なる場合もある(例えば、グループ 4 1 6 とサブグループ 4 2 8 および 4 3 2 との間の y - 軸に沿った距離は、グループ 4 2 0 とサブグループ 4 3 6 および 4 4 0 との間の y - 軸に沿った距離とは異なることも可能である)。y - 軸は、実施態様によっては、サブグループへの区分と関連付けられた t 1 スコアを表すこともできる。

【 0 0 7 5 】

[0081] 図 5 は、データを表示するユーザ・インターフェース 5 0 0 の一例である。ユーザ・インターフェース 5 0 0 は、第 1 表示部 5 0 4、第 2 表示部 5 0 8、および第 3 表示部 5 1 2 を含む。第 1 表示部 5 0 4 は、図 1 のグラフ 1 0 0 を含む。第 2 表示部 5 0 8 は、図 4 の分類ツリー 4 0 0 を含む。第 1 表示部 5 0 4 および第 2 表示部 5 0 8 は、ユーザが、分類ツリー 4 0 0 からのデータを、グラフ 1 0 0 上における対応するデータと素早くおよび/または視覚的に関連付けることを可能にする。例えば、ユーザが分類ツリー 4 0 0 のサブサブグループ 4 4 8 に含まれるデータに関心がある場合、ユーザは、第 2 表示部 4 0 0 において分類ツリー 4 0 0 からサブサブグループ 4 4 8 を選択することができる。選択されたサブサブグループ 4 4 8 は、第 2 表示部 5 0 8 において、楕円 5 1 6 またはその他のグラフィック指示手段によって強調される。対応する楕円 5 2 0、またはその他のグラフィック指示(例えば、強調)が、第 1 表示部 5 0 4 におけるグラフ 1 0 0 の上に現れる。楕円 5 2 0 またはその他のグラフィック指示は、第 1 サブサブグループ 4 4 8 において表されているデータを示す。実施形態によっては、サブサブグループ 4 4 8 の中にあるデータが密接にクラスタ化されていない場合や、楕円 5 2 0 またはその他のグラフィック指示によって識別が容易にできない場合もある。このような実施形態では、グラフ 1 0 0 におけるデータを異なる方法で(例えば、色または強調によって、あるいはラインまたは曲線の当てはめによって)表すことができる。

【 0 0 7 6 】

[0082] 同様に、図 4 の分類ツリー 4 0 0 のその他の階層レベルもユーザが選択し、第 1 表示部 5 0 4 に表示することができる。尚、追加の情報は、グラフ 1 0 0 上では第 1 表示部 5 0 4 において提示できることは言うまでもない。例えば、図示のように、グラフ 1 0 0 は、ライン 1 2 5 (第 1 階層レベル 4 0 4 における t 1 スコアに対応する)と、データ 1 0 5 を第 1 グループまたはセクション 1 3 0 および第 2 グループまたはセクション 1 3 5 に分離するライン 1 4 5 とを含む。第 1 グループ 1 3 0 は、第 1 階層レベル 4 0 4 における第 1 グループ 4 1 6 に対応し、第 2 グループ 1 3 5 は、第 1 階層レベル 4 0 4 における第 2 グループ 4 2 0 に対応する。第 1 グループ 4 1 6 を更に第 1 サブグループ 4 2 8 および第 2 サブグループ 4 3 2 に分類するとき、追加のライン(追加の t 1 スコアに対応するライン、およびそれを通過する垂直ライン)を第 1 表示部 5 0 4 におけるグラフ 1 0 0 に追加することができる。

【 0 0 7 7 】

[0083] 第 3 表示部 5 1 2 は、選択したデータを処理するためにプロセッサ(例えば、図 2 のプロセッサ 2 1 0)に命令を供給する、複数のユーザ選択可能ボタン 5 5 0 a ~ 5 5 0 h を含む。図示のように、サブサブグループ 4 4 8 が第 2 表示部 5 0 8 において選択されている。更に、ユーザは、ボタン 5 5 0 a ~ 5 5 0 h を通じて、サブサブグループ 4 4 8 の中にあるデータの特性を調査または評価することができる。例えば、ボタン 5 5 0 a は、サブサブグループ 4 4 8 内における t 1 値の分散を計算するモデルと関連付けられており、ユーザがボタン 5 5 0 a を選択したときに、t 1 値の分散が計算される。本明細書において言及する場合、分類ツリー内の特定のグループについての統計計算は、データ・マトリクスについて行われる計算を指し、その結果はマトリクスまたは 1 つの値(例えば、二乗の和)とすることができる。

【 0 0 7 8 】

[0084] ボタン 5 5 0 b は、サブサブグループ 4 4 8 内における t 1 値の標準偏差を計算

10

20

30

40

50

するモジュールと関連付けられており、ユーザがボタン 5 5 0 b を選択すると、 t_1 値の標準偏差が計算される。ボタン 5 5 0 c は、サブサブグループ 4 4 8 内における Y 値の分散（例えば、サブサブグループ 4 4 8 の第 2 データ・マトリクスまたは Y - マトリクスの分散）を計算するモジュールと関連付けられている。ボタン 5 5 0 c は、サブサブグループ 4 4 8 内にある異なる値（例えば、Y 値に対して t_1 値）が分析されることを除いて、ボタン 5 5 0 a と同様である。ボタン 5 5 0 d は、サブサブグループ 4 4 8 内における Y 値の標準偏差を計算するモジュールと関連付けられている。ボタン 5 5 0 d は、サブサブグループ 4 4 8 内における異なる値（例えば、Y 値に対して t_1 値）を分析することを除いて、ボタン 5 5 0 d と同様である。

【 0 0 7 9 】

[0085] ボタン 5 5 0 e および 5 5 0 f は、それぞれ、サブサブグループ 4 4 8 内における t_1 値および Y - 値の平均を計算するモジュールに関連付けられている。様々な計算技法のいずれかにしたがって、平均の代替物を計算することもでき、 t_1 値または Y - 値またはマトリクスの中央値またはモードを決定することを含む。実施形態によっては、計算技法はユーザが選択する。ユーザが選択したボタン 5 5 0 e ~ 5 5 0 f に応答して計算された t_1 値および / または Y - 値の平均は、更に別の評価（例えば、収量データ）または分析のためのデータを示すことができる。

【 0 0 8 0 】

[0086] ボタン 5 5 0 g は、 R^2 を計算するモジュールと関連付けられている。 R^2 は、Y - 値または収量値の変動を示し、PLS または OPLS 分析の正確さを判断するため、例えば、 t_1 スコア・ラインがどれ位正確に対応する Y - データに当てはまるかを判断するために用いることができる。実施形態によっては、 R^2 が多重相関係数として分かっている場合もある。ボタン 5 5 0 h は、 Q^2 を計算するモジュールと関連付けられている。 Q^2 は、クロス確認手順を用いる特定の PLS または OPLS モデルによって予測される、サブサブグループ 4 4 8（またはいずれかの Y - マトリクス）における全変動の端数 (fraction) を示す。

【 0 0 8 1 】

[0087] 実施形態によっては、ボタン 5 5 0 a ~ 5 5 0 h は、階層レベル内にある個々の分岐やグループではなく、階層レベル 4 0 4、4 1 2、4 2 4、および 4 4 4 と関連付けることができる場合もある。

【 0 0 8 2 】

[0088] 実施形態によっては、近似検索 (approximate search) を用いて最適な区分を判断する速度を増大させることができる。例えば、近似検索は、データの検索概要の多項式近似に基づくことができる。検索概要の近似の一例は、区分的二次多項式近似である。図 6 は、近似検索を用いるためのアルゴリズムの一例を示すフロー・チャート 6 0 0 である。

【 0 0 8 3 】

[0089] ステップ 6 0 4 において、各多項式近似に用いられる点 ($n_{p_{o1}}$) の数を選択する。点 ($n_{p_{o1}}$) の数は、ユーザが選択することができる。ユーザが点 ($n_{p_{o1}}$) の数に対して値を選択していない場合、点 ($n_{p_{o1}}$) の数のデフォルト値を用いる（ステップ 6 0 8）。例えば、点 ($n_{p_{o1}}$) の数のデフォルト値は、関数によって $n_{p_{o1}} = \min(11, \sqrt{N})$ と表すことができる。ここで、N はデータ集合における点の総数を表す。

【 0 0 8 4 】

[0090] 点 ($n_{p_{o1}}$) の数を決定した後、データを当てはめるために用いられる多項式区間 (polynomial piece) の数を決定する（ステップ 6 1 2）。多項式区間の数 $N_{p_{o1y}}$ を決定するために用いられる関数関係の一例は、 $N_{p_{o1y}} = \min(7, \text{整数}[2N/n_{p_{o1}}] - 1)$ である。

【 0 0 8 5 】

[0091] 多項式区間の数 $N_{p_{o1y}}$ を決定した後、近似の初期刻み長を計算する（ステップ 6 1 6）。刻み長は、最初の観察点および最後の観察点を除いた観察点（例えば、X - 変数）の範囲が包含されるように、そして各多項式区間が当該多項式区間の midpoint の各側において観察点の半分と重複するように決定される。

10

20

30

40

50

【0086】

[0092] 「切断表現」(cut expression)、即ち、式1におけるパラメータ「u」の値を、各多項式区間における点毎に計算する(ステップ620)。各多項式区間における点毎に式1を計算した後(ステップ620)、例えば、最小二乗当てはめ方法を用いて、二次多項式を各多項式区間に当てはめる(ステップ624)。ステップ624に関して、別の当てはめ技法を用いることもできる。当てはめられた多項式区間に対して最小値が生ずる「u」の値を計算する(ステップ628)。

【0087】

[0093] 最小値を計算した後(ステップ628)、刻み長を評価する(ステップ632)。当てはめプロセスの刻み長は、多項式における2点間の観察点の数を表す。刻み長が所定値(例えば、1)を超えていない場合、プロセスは終了する(ステップ636)。逆に、刻み長が所定値(例えば、1)を超過している場合、刻み長を短縮する(ステップ640)。例えば、刻み長を1/4に分割する(例えば、刻み長の値を4で除算する)ことができる。刻み長を短縮した後、新たな多項式区間を生成する(ステップ644)。新たな多項式区間は、パラメータ「u」の最小値を中心として位置付けられており、データのほぼ半分がこの多項式区間の中心のいずれの側にも来るようにしている。この新たな多項式区間について、この多項式区間における点毎に「u」の値を決定する(ステップ648)。ステップ648は、刻み長をステップ632において評価する前に行われるステップ620と動作が同様である。

10

【0088】

[0094] ステップ648においてパラメータ「u」の値を最小化した後、多項式を新たな多項式区間に当てはめ(ステップ652)、この新たな多項式区間におけるパラメータ「u」の最小値を計算する(ステップ656)。ステップ648の後、再度刻み長を評価し(ステップ632)、刻み長が所定数(例えば、1)を超えない場合、プロセスは終了する(ステップ636)。そうでない場合、刻み長が所定の閾値未満に減少するまで、プロセスを繰り返す。

20

【0089】

[0095] 以上で説明した技法は、デジタル電子回路において、またはコンピュータ・ハードウェア、ファームウェア、ソフトウェア、あるいはその組み合わせにおいて実現することができる。実施態様は、コンピュータ・プログラム・プロダクト、例えば、データ処理装置、例えば、プログラマブル・プロセッサ、コンピュータ、または多数のコンピュータによって実行するため、あるいはその動作を制御するために、情報担体、例えば、機械読み取り可能記憶デバイスに有体的に具現化されているコンピュータ・プログラムとすることができる。コンピュータ・プログラムは、コンパイル型言語またはインタプリタ型言語を含む、あらゆる形態のプログラミング言語でも書くことができ、単体プログラム、あるいはモジュール、コンポーネント、サブルーチン、あるいは計算環境において用いるのに適したその他のユニットを含む、あらゆる形態で展開することができる。コンピュータ・プログラムは、1カ所にある1つのコンピュータ、あるいは多数の箇所を跨って分散し通信ネットワークによって相互接続されている多数のコンピュータ上で実行するように展開することができる。

30

40

【0090】

[0096] 方法ステップは、1つ又は複数のプログラマブル・プロセッサによって実行することができる、入力データに作用し出力を発生することによって本技術の機能を実行するコンピュータ・プログラムを実行する。また、方法ステップは、特殊目的論理回路、例えば、FPGA(フィールド・プログラマブル・ゲート・アレイ)またはASIC(特定用途集積回路)によって実行することができる、そして装置はこの特殊目的論理回路として実現することができる。モジュールは、コンピュータ・プログラムの一部、および/またはその機能を実現するプロセッサ/特殊回路に言及することができる。

【0091】

[0097] コンピュータ・プログラムの実行に適したプロセッサは、一例として、汎用およ

50

び特殊目的用マイクロプロセッサの双方、ならびにあらゆる種類のデジタル・コンピュータの1つ又は複数のあらゆるプロセッサを含む。一般に、プロセッサは命令およびデータをリード・オンリ・メモリまたはランダム・アクセス・メモリあるいはその双方から受け取る。コンピュータの必須要素は、命令を実行するプロセッサと、命令およびデータを格納する1つ又は複数のメモリ・デバイスである。一般に、コンピュータは、データを格納する1つ又は複数の大容量記憶デバイス、例えば、磁気、光磁気ディスク、または光ディスクも含むか、またはこれらからデータを受け取りこれらにデータを転送する、あるいはこれらの双方を行うように動作的に結合される。データ送信および命令は、通信ネットワーク上においても生ずることができる。コンピュータ・プログラム命令およびデータを具現化するのに適した情報担体は、あらゆる形態の不揮発性メモリを含み、一例として、半導体メモリ・デバイス、例えば、EPROM、EEPROM、およびフラッシュ・メモリ・デバイス、磁気ディスク、例えば、内部ハード・ディスクまたはリムーバブル・ディスク、光磁気ディスク、ならびにCD-ROMおよびDVD-ROMディスクが含まれる。プロセッサおよびメモリは、特殊目的論理回路によって補強すること、または特殊目的論理回路の中に組み込むことができる。

10

20

30

40

50

【0092】

[0098] 「モジュール」および「機能」という用語は、本明細書において用いる場合、ある種のタスクを実行するソフトウェアまたはハードウェア・コンポーネントを意味するが、これに限定されるのではない。モジュールは、アドレス可能な記憶媒体上に常駐するように構成し、1つ又は複数のプロセッサにおいて実行するように構成することができる利点がある。モジュールは、汎用集積回路(「IC」)、FPGA、またはASICによって全体的または部分的に実現することができる。つまり、モジュールは、一例として、ソフトウェア・コンポーネント、オブジェクト指向ソフトウェア・コンポーネント、クラス・コンポーネントおよびタスク・コンポーネントというようなコンポーネント、プロセス、関数、属性、手順、サブルーチン、プログラム・コードのセグメント、ドライバ、ファームウェア、マイクロコード、回路、データ、データベース、データ構造、表、アレイ、ならびに変数を含むことができる。コンポーネントおよびモジュールにおいて設けられる機能は、もっと少ない数のコンポーネントおよびモジュールに組み合わせることができ、あるいは追加のコンポーネントおよびモジュールに更に分離することもできる。加えて、コンポーネントおよびモジュールは、多くの異なるプラットフォーム上で実現することができるという利点があり、これらのプラットフォームには、コンピュータ、コンピュータ・サーバ、アプリケーション対応(application-enabled)スイッチまたはルータのようなデータ通信インフラストラクチャ機器、あるいは公衆または個人電話スイッチまたは個人分岐交換機(「PBX」(private branch exchange))のような電気通信インフラストラクチャ機器が含まれる。これらの場合のいずれにおいても、選択したプラットフォームにネイティブなアプリケーションを書くことによって、あるいはプラットフォームを1つ又は複数の外部アプリケーション・エンジンにインターフェースすることのいずれかによって、実現を達成することができる。

【0093】

[0099] ユーザとの双方向処理に備えるために、前述した技法は、ユーザに情報を表示するディスプレイ・デバイス、例えば、CRT(陰極線管)またはLCD(液晶ディスプレイ)モニター、ならびにキーボードおよびポインティング・デバイス、例えば、マウスまたはトラックボールを有するコンピュータ上で実現することができる。ユーザは、ポインティング・デバイスによってコンピュータに入力を供給することができる(例えば、ユーザ・インターフェース・エレメントと双方向処理を行う)。ユーザとの双方向処理に備えるためには、他の種類のデバイスも用いることができる。例えば、ユーザに宛てるフィードバックは、あらゆる形態の感覚的フィードバック、例えば、視覚フィードバック、聴覚フィードバック、または触覚フィードバックとすることができ、ユーザからの入力は、音響入力、音声入力、または接触入力を含む、あらゆる形態で受け取ることができる。

【0094】

[0100] 前述した技法は、分散型計算システムにおいて実現することができる。分散型計算システムは、例えば、データ・サーバのようなバック・エンド・コンポーネント、および/またはミドルウェア・コンポーネント、例えば、アプリケーション・サーバ、および/またはフロント・エンド・コンポーネント、例えば、グラフィカル・ユーザ・インターフェースおよび/またはウェブ・ブラウザを有し、これを通じてユーザが実施態様例と双方向処理することができるクライアント・コンピュータ、あるいは、このようなバック・エンド、ミドルウェア、またはフロント・エンド・コンポーネントのあらゆる組み合わせを含む。本システムのコンポーネントは、あらゆる形態または媒体のデジタル・データ通信、例えば、通信ネットワークによって相互接続することができる。通信チャネルとも呼ばれる通信ネットワークの例には、ローカル・エリア・ネットワーク(「LAN」)およびワイド・エリア・ネットワーク(「WAN」)、例えば、インターネットが含まれ、そして有線ネットワークおよびワイヤレス・ネットワークの双方が含まれる。例の中には、通信ネットワークが仮想ローカル・エリア・ネットワーク(「VLAN」)のような仮想ネットワークまたはサブネットワークを特徴とすることができるものもある。特に明確に示さない限り、通信ネットワークは、PSTNの全部または一部、例えば、特定の電気通信事業者が所有する一部も含むことができる。

10

【0095】

[0101] 前述の計算システムは、クライアントおよびサーバを含むことができる。クライアントおよびサーバは、一般に、互いに離れており、通信ネットワークを通じて双方向処理するのが通例である。クライアントおよびサーバの関係は、それぞれのコンピュータ上において実行し、互いにクライアント・サーバ関係を有するコンピュータ・プログラムによって生ずる。

20

【0096】

[0102] 種々の実施形態は、交信状態にあるものまたは1系統以上の通信経路によって接続されたものとして図示されている。通信経路は、特定のデータ転送媒体に限定されることはない。情報は、電気信号、光信号、音響信号、物理信号、熱信号、またはそのあらゆる組み合わせを用いて、通信経路上で送信することができる。通信経路は、多数の通信チャネル、例えば、データ・フローの容量が同一であるまたは異なる、多重化チャネルを含むことができる。

【0097】

[0103] 図示したユーザ・インターフェース機構(user interface feature)のパラメータを設定するためには、多数のユーザ入力を用いることができる。このような入力の例には、ボタン、ラジオ・ボタン、アイコン、チェック・ボックス、コンボ・ボックス、メニュー、テキスト・ボックス、ツールチップ、トグル・スイッチ、ボタン、スクロール・バー、ツールバー、ステータス・バー、ウィンドウ、あるいはユーザが、本明細書に記載したモジュールまたはシステムのいずれとでも通信すること、および/またはデータを提供することを可能にするユーザ・インターフェースと関連付けられている、その他の適したアイコンまたはウィジェット(widget)が含まれる。

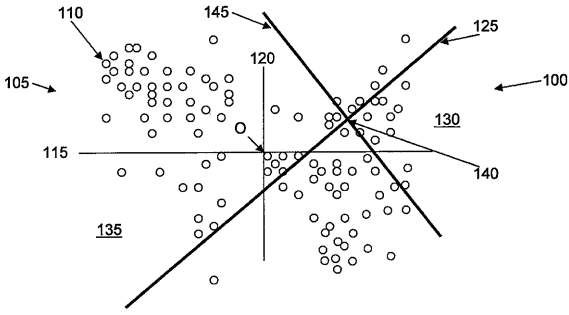
30

【0098】

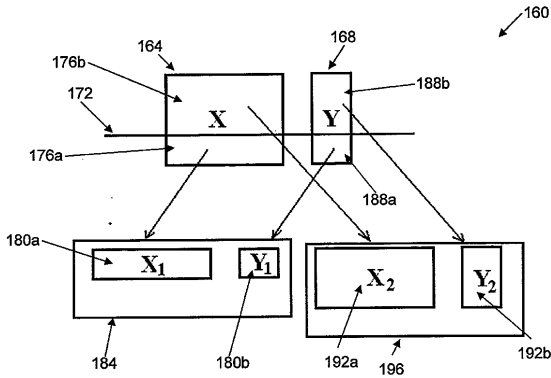
[0104] 以上、具体的な実施形態を参照しながら、本発明について特定的に示し説明したが、添付した特許請求の範囲によって定められる発明の主旨および範囲から逸脱することなく、本発明の形態および詳細には種々の変更が可能であることは、当業者には理解されてしかるべきである。

40

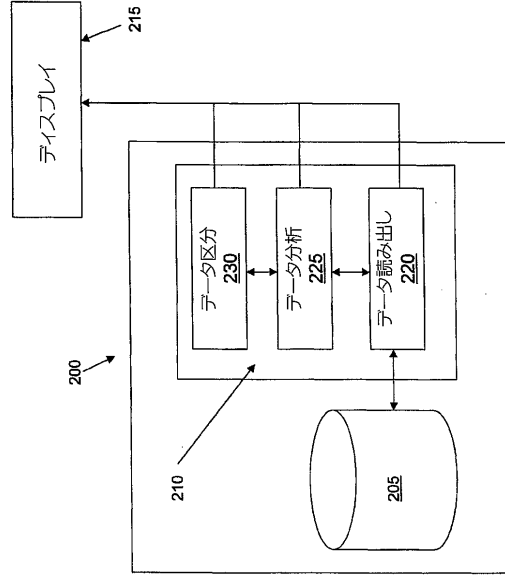
【図1A】



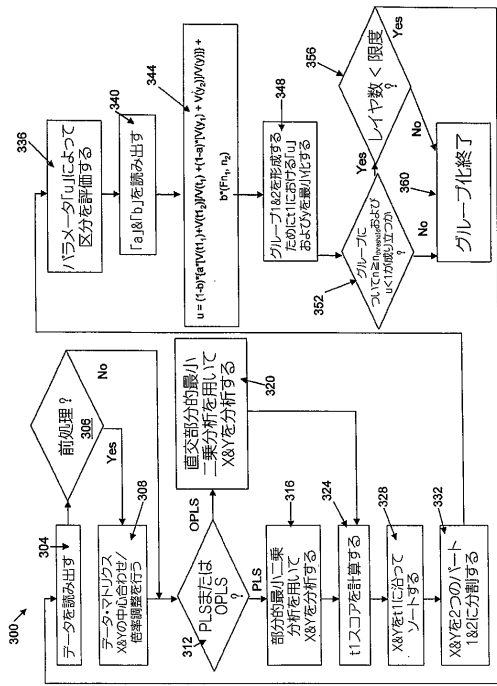
【図1B】



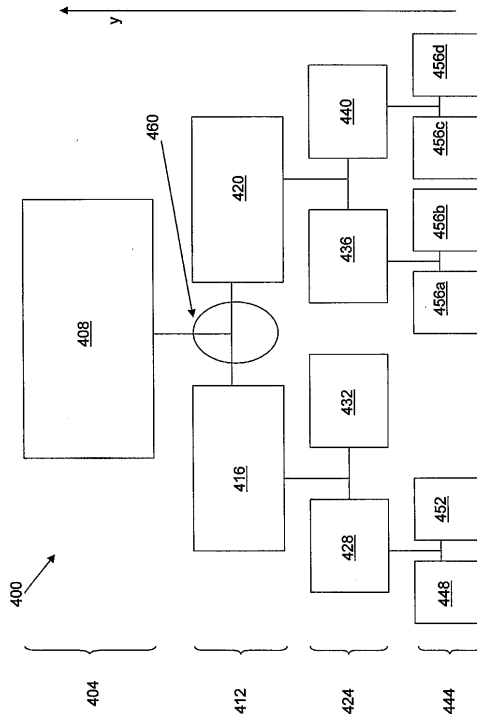
【図2】



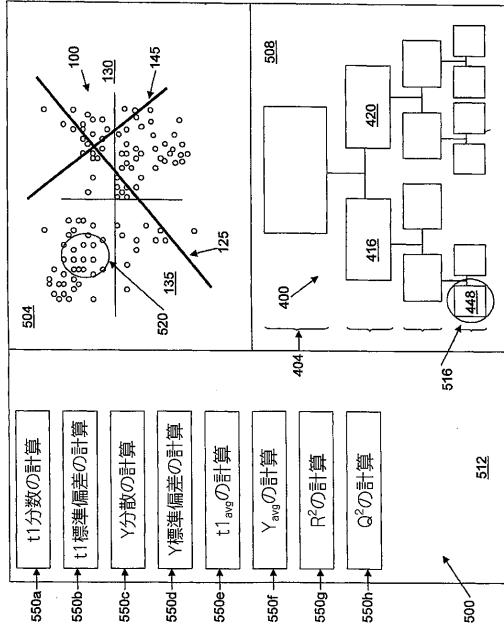
【図3】



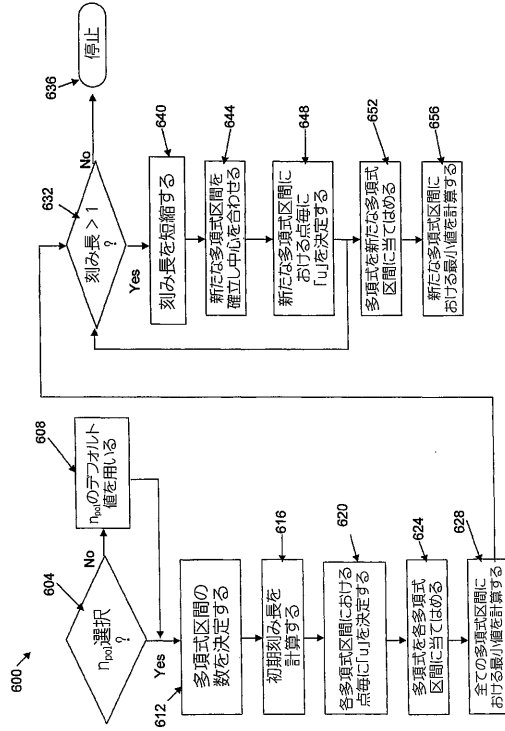
【図4】



【図 5】



【図 6】



フロントページの続き

(74)代理人 100118902

弁理士 山本 修

(74)代理人 100162846

弁理士 大牧 綾子

(72)発明者 ウォルド, スヴァント・ブジャルヌ

アメリカ合衆国ニューハンプシャー州03049, ホリス, パイン・ヒル・ロード 42

(72)発明者 トリゲ, ユーハン

スウェーデン国 エスイー - 904 35, スウェーデン, ウメア, オクスベルスヴェーゲン 20

(72)発明者 エリクソン, レンナート

スウェーデン国 エスイー - 90751, スウェーデン, ウメア, スヨフルヴァーゲン 85