

(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2020-161147

(P2020-161147A)

(43) 公開日 令和2年10月1日(2020.10.1)

(51) Int.Cl.
G06F 16/21 (2019.01)

F I
G06F 16/21

テーマコード (参考)
5B175

審査請求 有 請求項の数 40 O L 外国語出願 (全 40 頁)

(21) 出願番号 特願2020-88498 (P2020-88498)
(22) 出願日 令和2年5月20日 (2020.5.20)
(62) 分割の表示 特願2017-559576 (P2017-559576) の分割
原出願日 平成28年6月10日 (2016.6.10)
(31) 優先権主張番号 62/174,997
(32) 優先日 平成27年6月12日 (2015.6.12)
(33) 優先権主張国・地域又は機関 米国 (US)
(31) 優先権主張番号 15/175,793
(32) 優先日 平成28年6月7日 (2016.6.7)
(33) 優先権主張国・地域又は機関 米国 (US)

(71) 出願人 509123208
アビニシオ テクノロジー エルエルシー
アメリカ合衆国 02421 マサチュー
セッツ州 レキシントン スプリング ス
トリート 201
(74) 代理人 100079108
弁理士 稲葉 良幸
(74) 代理人 100109346
弁理士 大貫 敏史
(74) 代理人 100117189
弁理士 江口 昭彦
(74) 代理人 100134120
弁理士 内藤 和彦

(特許庁注：以下のものは登録商標)

最終頁に続く

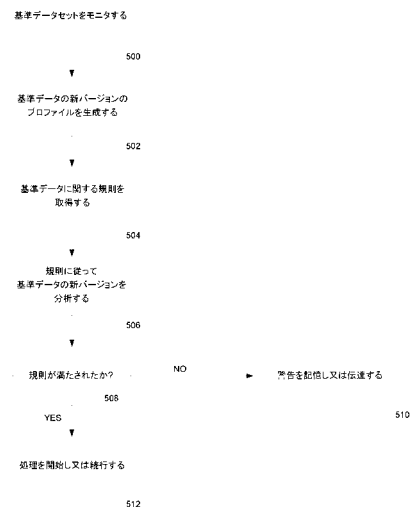
(54) 【発明の名称】 データ品質分析

(57) 【要約】 (修正有)

【課題】データ品質分析の向上を図る方法、非一時的コンピュータ可読媒体及び計算システムを提供する。

【解決手段】方法は、データセットを分析することであって、データセットのフィールドについての基準プロファイルを決定することと、決定した基準プロファイルに基づいて、データセットのフィールドについてのデータ品質規則を決定することと、を含む。データセットのフィールドについてのデータ品質規則は、フィールドについての基準プロファイルと、データセットの一つ以上のデータ記録のフィールドについてのプロファイルとの間の許容偏差、データセットのデータ記録のフィールドのデータ要素についての許容値又はデータセットのデータ記録のフィールドのデータ要素についての禁止値、の一つ以上を示す。

【選択図】 図 1 2



【特許請求の範囲】

【請求項 1】

データ処理システムによって生成される出力データセットを示す情報を受信することと

、
前記出力データセットに関係するデータ系列情報に基づき、前記出力データセットが依拠する1つ又は複数のアップストリームデータセットを識別することと、

前記出力データセットが依拠する前記識別された1つ又は複数のアップストリームデータセットの1つ又は複数进行分析することであって、前記1つ又は複数のアップストリームデータセットのうちの特定のアップストリームデータセットごとに、

(i) 前記特定のアップストリームデータセットのプロファイルと前記特定のアップストリームデータセットに関する基準プロファイルとの間の許容偏差を示す第1の規則、及び

(i i) 前記特定のアップストリームデータセット内の1つ又は複数のデータ要素のそれぞれに関する1つ又は複数の許容値又は禁止値を示す第2の規則

のうちの1つ又は複数適用し、前記1つ又は複数の規則を適用した結果に基づき、前記アップストリームデータセットの1つ又は複数を選択すること、

を含む、分析することと、

前記選択された1つ又は複数のアップストリームデータセットに関連する情報を出力することと、

を含む、方法。

【請求項 2】

前記第1の規則及び前記第2の規則の1つ又は複数が自動で生成される、請求項1に記載の方法。

【請求項 3】

前記第1の規則が、前記特定のアップストリームデータセットの履歴プロファイルの自動分析に基づいて自動で生成される、請求項2に記載の方法。

【請求項 4】

前記基準プロファイルが、前記特定のアップストリームデータセットに関する履歴平均プロファイルに基づく、請求項3に記載の方法。

【請求項 5】

前記第2の規則が、前記特定のアップストリームデータセット内の前記1つ又は複数のデータ要素に関する履歴値の自動分析に基づいて自動で生成される、請求項2に記載の方法。

【請求項 6】

前記許容値又は禁止値が前記自動分析に基づいて決定される、請求項5に記載の方法。

【請求項 7】

前記第1の規則及び前記第2の規則の1つ又は複数がユーザによって指定される、請求項1に記載の方法。

【請求項 8】

前記第1の規則及び前記第2の規則の1つ又は複数の指定を、ユーザインタフェースを通して受信することを含む、請求項1に記載の方法。

【請求項 9】

データ系列情報が、前記出力データセットが依拠する1つ又は複数のデータセット、前記出力データセットに依拠する1つ又は複数のデータセット、又はその両方を示す、請求項1に記載の方法。

【請求項 10】

前記データセットの部分集合を識別するために前記1つ又は複数のデータセットのそれぞれ进行分析することが、前記1つ又は複数のデータセットの何れが誤り又は起こり得る誤りを有するかを判定することを含み、

前記方法が、前記部分集合に関して誤り又は起こり得る誤りを有する前記データセット

10

20

30

40

50

を選択することを含む、請求項 1 に記載の方法。

【請求項 1 1】

前記データセットの部分集合を識別するために前記 1 つ又は複数のデータセットのそれぞれを分析することが、特定のデータセットであって、前記特定のデータセットの前記プロファイルと前記特定のデータセットに関する前記基準プロファイルとの間の偏差が、前記対応する第 1 の規則によって示される前記許容偏差を上回る、特定のデータセットを識別することを含む、

前記方法が、前記部分集合のために前記特定のデータセットを選択することを含む、請求項 1 に記載の方法。

【請求項 1 2】

前記データセットの部分集合を識別するために前記 1 つ又は複数のデータセットのそれぞれを分析することが、前記対応する第 2 の規則によって示される前記許容値又は禁止値を満たさない値を有するデータ要素を有する特定のデータセットを識別することを含む、

前記方法が、前記部分集合のために前記特定のデータセットを選択することを含む、請求項 1 に記載の方法。

【請求項 1 3】

前記出力データセット内のデータ要素を識別することを含む、前記出力データセットが依拠する前記 1 つ又は複数のデータセットを識別することが、前記出力データセット内の前記識別されたデータ要素に影響を及ぼすデータセットを識別することを含む、請求項 1 に記載の方法。

【請求項 1 4】

前記出力データセット内のデータ要素を識別することが、誤り又は起こり得る誤りを有するデータ要素を識別することを含む、請求項 1 3 に記載の方法。

【請求項 1 5】

前記アップストリームデータセットの 1 つ又は複数のプロファイルを生成することを含む、請求項 1 に記載の方法。

【請求項 1 6】

特定のデータセットのプロファイルを生成することが、前記特定のデータセットの新バージョンが受信されるときに前記特定のデータセットの新規プロファイルを生成することを含む、請求項 1 5 に記載の方法。

【請求項 1 7】

特定のデータセットに関する前記基準プロファイルが、前記特定のデータセットの 1 つ又は複数の過去のプロファイルから導出される、請求項 1 に記載の方法。

【請求項 1 8】

前記データセットの部分集合に関連する情報を出力することが、前記部分集合の前記データセットのそれぞれの識別子を出力することを含む、請求項 1 に記載の方法。

【請求項 1 9】

前記データセットの部分集合に関連する情報を出力することが、前記部分集合の前記データセットのそれぞれに関連する誤り又は起こり得る誤りの標識を出力することを含む、請求項 1 に記載の方法。

【請求項 2 0】

前記データ処理システムの表現をユーザインタフェース上で表示することを含む、前記データセットの部分集合に関連する情報を出力することが、前記データセットの部分集合の特定のデータセットの表現の近くに前記部分集合の前記特定のデータセットに関連する情報を表示することを含む、請求項 1 に記載の方法。

【請求項 2 1】

前記部分集合の前記特定のデータセットに関連する前記表示された情報が、前記特定のデータセットの前記プロファイルと前記特定のデータセットに関する前記基準プロファイルとの間の偏差を示す値を含む、請求項 2 0 に記載の方法。

【請求項 2 2】

10

20

30

40

50

前記部分集合の前記特定のデータセットに関連する前記表示された情報が、前記対応する第2の規則によって示される前記許容値又は禁止値を満たさない前記特定のデータセット内のデータ要素の数を表す値を含む、請求項20に記載の方法。

【請求項23】

前記データセットの部分集合に関する情報を示す情報バブル又はポップアップウィンドウを表示することを含む、請求項20に記載の方法。

【請求項24】

ユーザが規則を追加するか、規則を修正するか、又は規則を除去することを可能にするためのユーザインタフェースを提供することを含む、請求項1に記載の方法。

【請求項25】

前記データセットが1つ又は複数のソースデータセット及び1つ又は複数の基準データセットを含み、前記ソースデータセットが、前記データ処理システムによって処理されるデータ要素を含み、前記基準データセットが、前記ソースデータセット内の前記データ要素を処理する際に前記データ処理システムによって参照される基準値を含む、請求項1に記載の方法。

【請求項26】

前記基準データセットが、前記データ処理システムに関連する企業体に関連するデータを含み、及び前記ソースデータセットが、前記企業体の顧客に関連するデータを含む、請求項25に記載の方法。

【請求項27】

前記データ処理システムが変換要素を含み、及び前記方法が、前記出力データセットに影響を及ぼす1つ又は複数の変換要素を前記データ系列情報に基づいて識別することを含む、請求項1に記載の方法。

【請求項28】

前記変換要素の何れの1つ又は複数が誤り又は起こり得る誤りを有するかを判定することを含む、請求項27に記載の方法。

【請求項29】

特定のデータ処理要素が誤り又は起こり得る誤りを有するかどうかを、前記特定の変換要素に関連する実装日に基づいて判定することを含む、請求項28に記載の方法。

【請求項30】

データ処理システムによって生成される出力データセットを示す情報を受信することと、

前記出力データセットに関係するデータ系列情報に基づき、前記出力データセットが依拠する1つ又は複数のアップストリームデータセットを識別することと、

前記出力データセットが依拠する前記識別された1つ又は複数のアップストリームデータセットの1つ又は複数を選択することであって、前記1つ又は複数のアップストリームデータセットのうちの前記特定のアップストリームデータセットごとに、

(i) 前記特定のアップストリームデータセットのプロファイルと前記特定のアップストリームデータセットに関する基準プロファイルとの間の許容偏差を示す第1の規則、及び

(ii) 前記特定のアップストリームデータセット内の1つ又は複数のデータ要素のそれぞれに関する1つ又は複数の許容値又は禁止値を示す第2の規則

のうち1つ又は複数を選択し、前記1つ又は複数の規則を適用した結果に基づき、前記アップストリームデータセットの1つ又は複数を選択すること、

を含む、分析することと、

前記選択された1つ又は複数のアップストリームデータセットに関連する情報を出力することと、

を計算システムに行わせるための命令を記憶する、非一時的コンピュータ可読媒体。

【請求項31】

メモリに結合されるプロセッサを含む計算システムであって、前記プロセッサ及びメモ

10

20

30

40

50

りは、

データ処理システムによって生成される出力データセットを示す情報を受信することと

、
前記出力データセットに関係するデータ系列情報に基づき、前記出力データセットが依拠する1つ又は複数のアップストリームデータセットを識別することと、

前記出力データセットが依拠する前記識別された1つ又は複数のアップストリームデータセットの1つ又は複数进行分析することであって、前記1つ又は複数のアップストリームデータセットのうち特定のアップストリームデータセットごとに、

(i) 前記特定のアップストリームデータセットのプロファイルと前記特定のアップストリームデータセットに関する基準プロファイルとの間の許容偏差を示す第1の規則、及び

10

(ii) 前記特定のアップストリームデータセット内の1つ又は複数のデータ要素のそれぞれに関する1つ又は複数の許容値又は禁止値を示す第2の規則

のうちの1つ又は複数適用し、前記1つ又は複数の規則を適用した結果に基づき、前記アップストリームデータセットの1つ又は複数を選択すること、

を含む、分析することと、

前記選択された1つ又は複数のアップストリームデータセットに関連する情報を出力することと、

を行うように構成される、計算システム。

【請求項32】

20

データ処理システムによって生成される出力データセットを示す情報を受信するための手段と、

前記出力データセットに関係するデータ系列情報に基づき、前記出力データセットが依拠する1つ又は複数のアップストリームデータセットを識別するための手段と、

前記出力データセットが依拠する前記識別された1つ又は複数のアップストリームデータセットの1つ又は複数进行分析するための手段であって、前記分析することが、前記1つ又は複数のアップストリームデータセットのうち特定のアップストリームデータセットごとに、

(i) 前記特定のアップストリームデータセットのプロファイルと前記特定のアップストリームデータセットに関する基準プロファイルとの間の許容偏差を示す第1の規則、及び

30

(ii) 前記特定のアップストリームデータセット内の1つ又は複数のデータ要素のそれぞれに関する1つ又は複数の許容値又は禁止値を示す第2の規則

のうちの1つ又は複数適用し、前記1つ又は複数の規則を適用した結果に基づき、前記アップストリームデータセットの1つ又は複数を選択すること、

を含む、手段と、

前記選択された1つ又は複数のアップストリームデータセットに関連する情報を出力するための手段と、

を含む、計算システム。

【請求項33】

40

データ処理システムのダウンストリームデータセットのデータ要素内の誤り又は起こり得る誤りを識別すると、前記ダウンストリームデータセットに関係するデータ系列情報に基づき、前記データ要素に影響を及ぼす1つ又は複数のアップストリームデータセットを自動で識別することと、

前記識別されたアップストリームデータセットのそれぞれの現在のプロファイル及び基準プロファイル进行分析することを含む、前記アップストリームデータセットの何れが誤り又は起こり得る誤りを有するかを判定することと、

誤りを有するか又は誤りを有する可能性が高いと判定される前記アップストリームデータセットのそれぞれに関連する情報を出力することと、

を含む、方法。

50

【発明の詳細な説明】

【技術分野】

【0001】

背景

本明細書は、データ品質分析に関する。データセットのデータ品質は、データセット内のデータ記録が誤りを有するかどうかの指標である。多くの場合、データセットの処理中に誤りが生じる場合、そのデータセットのデータ品質は低い。

【発明の概要】

【課題を解決するための手段】

【0002】

要約

一般的な態様では、方法は、データ処理システムによって生成される出力データセットを示す情報を受信することと、出力データセットに係るデータ系列情報に基づき、出力データセットが依拠する1つ又は複数のアップストリームデータセットを識別することと、出力データセットが依拠する識別された1つ又は複数のアップストリームデータセットの1つ又は複数を含む。分析することは、1つ又は複数のアップストリームデータセットのうち特定のアップストリームデータセットごとに、(i)特定のアップストリームデータセットのプロファイルと特定のアップストリームデータセットに関する基準プロファイルとの間の許容偏差を示す第1の規則、及び(ii)特定のアップストリームデータセット内の1つ又は複数のデータ要素のそれぞれに関する1つ又は複数の許容値又は禁止値を示す第2の規則のうち1つ又は複数を選択することを含む。この方法は、選択された1つ又は複数のアップストリームデータセットに関連する情報を出力することを含む。

【0003】

実施形態は、以下の特徴の1つ又は複数を含み得る。

【0004】

第1の規則及び第2の規則の1つ又は複数が自動で生成される。第1の規則は、特定のアップストリームデータセットの履歴プロファイルの自動分析に基づいて自動で生成される。基準プロファイルは、特定のアップストリームデータセットに関する履歴平均プロファイルに基づく。第2の規則は、特定のアップストリームデータセット内の1つ又は複数のデータ要素に関する履歴値の自動分析に基づいて自動で生成される。許容値又は禁止値は自動分析に基づいて決定される。

【0005】

第1の規則及び第2の規則の1つ又は複数がユーザによって指定される。

【0006】

この方法は、第1の規則及び第2の規則の1つ又は複数の指定を、ユーザインタフェースを通して受信することを含む。

【0007】

データ系列情報は、出力データセットが依拠する1つ又は複数のデータセット、出力データセットに依拠する1つ又は複数のデータセット、又はその両方を示す。

【0008】

データセットの部分集合を識別するために1つ又は複数のデータセットのそれぞれを分析することは、1つ又は複数のデータセットの何れが誤り又は起こり得る誤りを有するかを判定することを含み、この方法は、部分集合に関して誤り又は起こり得る誤りを有するデータセットを選択することを含む。

【0009】

データセットの部分集合を識別するために1つ又は複数のデータセットのそれぞれを分析することは、特定のデータセットであって、特定のデータセットのプロファイルと特定のデータセットに関する基準プロファイルとの間の偏差が、対応する第1の規則によって

10

20

30

40

50

示される許容偏差を上回る、特定のデータセットを識別することを含み、この方法は、部分集合のために特定のデータセットを選択することを含む。

【0010】

データセットの部分集合を識別するために1つ又は複数のデータセットのそれぞれを分析することは、対応する第2の規則によって示される許容値又は禁止値を満たさない値を有するデータ要素を有する特定のデータセットを識別することを含み、この方法は、部分集合のために特定のデータセットを選択することを含む。

【0011】

この方法は、出力データセット内のデータ要素を識別することを含み、出力データセットが依拠する1つ又は複数のデータセットを識別することは、出力データセット内の識別されたデータ要素に影響を及ぼすデータセットを識別することを含む。出力データセット内のデータ要素を識別することは、誤り又は起こり得る誤りを有するデータ要素を識別することを含む。

10

【0012】

この方法は、アップストリームデータセットの1つ又は複数のプロファイルを生成することを含む。特定のデータセットのプロファイルを生成することは、特定のデータセットの新バージョンが受信されるときに特定のデータセットの新規プロファイルを生成することを含む。

【0013】

特定のデータセットに関する基準プロファイルは、特定のデータセットの1つ又は複数の過去のプロファイルから導出される。

20

【0014】

データセットの部分集合に関連する情報を出力することは、部分集合のデータセットのそれぞれの識別子を出力することを含む。

【0015】

データセットの部分集合に関連する情報を出力することは、部分集合のデータセットのそれぞれに関連する誤り又は起こり得る誤りの標識を出力することを含む。

【0016】

この方法は、データ処理システムの表現をユーザインタフェース上で表示することを含み、データセットの部分集合に関連する情報を出力することは、データセットの部分集合の特定のデータセットの表現の近くに部分集合の特定のデータセットに関連する情報を表示することを含む。部分集合の特定のデータセットに関連する表示された情報は、特定のデータセットのプロファイルと特定のデータセットに関する基準プロファイルとの間の偏差を示す値を含む。部分集合の特定のデータセットに関連する表示された情報は、対応する第2の規則によって示される許容値又は禁止値を満たさない特定のデータセット内のデータ要素の数を表す値を含む。この方法は、データセットの部分集合に関する情報を示す情報バブル又はポップアップウィンドウを表示することを含む。

30

【0017】

この方法は、ユーザが規則を追加するか、規則を修正するか、又は規則を除去することを可能にするためのユーザインタフェースを提供することを含む。

40

【0018】

データセットは1つ又は複数のソースデータセット及び1つ又は複数の基準データセットを含み、ソースデータセットは、データ処理システムによって処理されるデータ要素を含み、基準データセットは、ソースデータセット内のデータ要素を処理する際にデータ処理システムによって参照される基準値を含む。基準データセットは、データ処理システムに関連する企業体に関連するデータを含み、及びソースデータセットは、企業体の顧客に関連するデータを含む。

【0019】

データ処理システムは変換要素を含み、及びこの方法は、出力データセットに影響を及ぼす1つ又は複数の変換要素をデータ系列情報に基づいて識別することを含む。この方法

50

は、変換要素の何れの1つ又は複数が誤り又は起こり得る誤りを有するかを判定することを含む。この方法は、特定のデータ処理要素が誤り又は起こり得る誤りを有するかどうかを、特定の変換要素に関連する実装日に基づいて判定することを含む。

【0020】

一般的な態様では、非一時的コンピュータ可読媒体は、データ処理システムによって生成される出力データセットを示す情報を受信することと、出力データセットに関係するデータ系列情報に基づき、出力データセットが依拠する1つ又は複数のアップストリームデータセットを識別することと、出力データセットが依拠する識別された1つ又は複数のアップストリームデータセットの1つ又は複数を分析することとを計算システムに行わせるための命令を記憶する。分析することは、1つ又は複数のアップストリームデータセットのうち特定のアップストリームデータセットごとに、(i)特定のアップストリームデータセットのプロファイルと特定のアップストリームデータセットに関する基準プロファイルとの間の許容偏差を示す第1の規則、及び(ii)特定のアップストリームデータセット内の1つ又は複数のデータ要素のそれぞれに関する1つ又は複数の許容値又は禁止値を示す第2の規則のうち1つ又は複数を適用し、且つ1つ又は複数の規則を適用した結果に基づき、アップストリームデータセットの1つ又は複数を選択することを含む。命令は、選択された1つ又は複数のアップストリームデータセットに関連する情報を計算システムに出力させる。

10

【0021】

一般的な態様では、計算システムは、メモリに結合されるプロセッサを含む。プロセッサ及びメモリは、データ処理システムによって生成される出力データセットを示す情報を受信することと、出力データセットに関係するデータ系列情報に基づき、出力データセットが依拠する1つ又は複数のアップストリームデータセットを識別することと、出力データセットが依拠する識別された1つ又は複数のアップストリームデータセットの1つ又は複数を分析することとを行うように構成される。分析することは、1つ又は複数のアップストリームデータセットのうち特定のアップストリームデータセットごとに、(i)特定のアップストリームデータセットのプロファイルと特定のアップストリームデータセットに関する基準プロファイルとの間の許容偏差を示す第1の規則、及び(ii)特定のアップストリームデータセット内の1つ又は複数のデータ要素のそれぞれに関する1つ又は複数の許容値又は禁止値を示す第2の規則のうち1つ又は複数を適用し、且つ1つ又は複数の規則を適用した結果に基づき、アップストリームデータセットの1つ又は複数を選択することを含む。プロセッサ及びメモリは、選択された1つ又は複数のアップストリームデータセットに関連する情報を出力するように構成される。

20

30

【0022】

一般的な態様では、計算システムは、データ処理システムによって生成される出力データセットを示す情報を受信するための手段と、出力データセットに関係するデータ系列情報に基づき、出力データセットが依拠する1つ又は複数のアップストリームデータセットを識別するための手段と、出力データセットが依拠する識別された1つ又は複数のアップストリームデータセットの1つ又は複数を分析するための手段とを含む。分析することは、1つ又は複数のアップストリームデータセットのうち特定のアップストリームデータセットごとに、(i)特定のアップストリームデータセットのプロファイルと特定のアップストリームデータセットに関する基準プロファイルとの間の許容偏差を示す第1の規則、及び(ii)特定のアップストリームデータセット内の1つ又は複数のデータ要素のそれぞれに関する1つ又は複数の許容値又は禁止値を示す第2の規則のうち1つ又は複数を適用し、且つ1つ又は複数の規則を適用した結果に基づき、アップストリームデータセットの1つ又は複数を選択することを含む。この計算システムは、選択された1つ又は複数のアップストリームデータセットに関連する情報を出力するための手段を含む。

40

【0023】

一般的な態様では、方法は、データ処理システムのダウンストリームデータセットのデータ要素内の誤り又は起こり得る誤りを識別すると、ダウンストリームデータセットに関

50

係するデータ系列情報に基づき、データ要素に影響を及ぼす1つ又は複数のアップストリームデータセットを自動で識別することと、識別されたアップストリームデータセットのそれぞれの現在のプロファイル及び基準プロファイルを分析することを含む、アップストリームデータセットの何れが誤り又は起こり得る誤りを有するかを判定することと、誤りを有するか又は誤りを有する可能性が高いと判定されるアップストリームデータセットのそれぞれに関連する情報を出力することを含む。

【0024】

態様は以下の利点の1つ又は複数を含み得る。

【0025】

本明細書に記載の手法は、データ分析官又はアプリケーション開発者等のユーザがデータ品質問題の根本的原因を迅速に識別することを補助し得る。例えば、データ処理システム内の基準データは頻繁に更新されるが、必ずしも導入前に完全に検査されない可能性がある。基準データ内の誤りは、基準データを使用して処理されるダウストリームデータ内のデータ品質問題を引き起こし得る。ダウストリームデータセット内のデータ品質問題の根本的原因を分析することは、ダウストリームデータセットのデータ品質に影響を及ぼしている可能性があるデータ品質問題を有する基準データ又は他のアップストリームデータを識別するのを補助し得る。潜在的なデータ品質問題をユーザに通知することは、ユーザがデータ処理をプロアクティブに管理することを補助し得る。

10

【0026】

本発明の他の特徴及び利点が以下の説明及び特許請求の範囲から明らかになる。

20

【図面の簡単な説明】

【0027】

【図1】データ系列図である。

【図2】データ系列図である。

【図3A】データ系列図である。

【図3B】データ系列図である。

【図4】ユーザインタフェースの図である。

【図5】システム図である。

【図6】ユーザインタフェースの図である。

【図7】データ処理システムの図である。

30

【図8A】データ処理システムの図である。

【図8B】データ処理システムの図である。

【図8C】記録の一例である。

【図9A】データ処理システムの図である。

【図9B】データ処理システムの図である。

【図10A】データ処理システムの図である。

【図10B】記録の一例である。

【図11】フローチャートである。

【図12】フローチャートである。

【図13】フローチャートである。

40

【図14】フローチャートである。

【図15】フローチャートである。

【図16】システム図である。

【発明を実施するための形態】

【0028】

説明

データ系列の分析に基づいてデータ品質問題の根本的原因を識別するための手法をここで説明する。データ品質問題がダウストリームデータセット内で識別される場合、ダウストリームデータセットの導出元であるアップストリームデータセット及びアップストリーム変換要素（アップストリームデータ系列要素と呼ばれる場合もある）が識別される

50

。ダウンストリームデータセット内のデータ品質問題に寄与したデータ品質問題をそれ自体が有し得るアップストリームデータ系列要素の1つ又は複数を識別するために、各アップストリームデータ系列要素の品質が評価される。一部の例では、データセットがデータ品質問題を有するかどうかを判定するために、各アップストリームデータセットを特徴付けるプロファイルがそのデータセットに関する履歴平均プロファイル等の基準プロファイルと比較される。一部の例では、データセットがデータ品質問題を有するかどうかを判定するために、アップストリームデータセットのフィールド内の値がそのフィールドに関する1つ又は複数の許容値又は禁止値と比較される。

【0029】

データ系列とは、データ処理システムによって処理されるデータ記録のライフサイクルを記述する情報である。所与のデータセットのためのデータ系列情報は、所与のデータセットが依拠する1つ又は複数のアップストリームデータセットの識別子、所与のデータセットに依拠する1つ又は複数のダウンストリームデータセット、及びデータを処理して所与のデータセットを生成する1つ又は複数の変換を含む。アップストリームデータセットに依拠するダウンストリームデータセットとは、データ処理システムによるアップストリームデータセットの処理がダウンストリームデータセットの生成を直接又は間接的にもたらすことを意味する。生成されるダウンストリームデータセットは、データ処理システムから出力されるデータセット（出力データセットと呼ばれる場合もある）とすることができ、又はデータ処理システムによって更に処理されるデータセット（中間データセットと呼ばれる場合もある）とすることができる。アップストリームデータセットは、データ処理システム内に入力されるデータセット（入力データセット又は基準データセットと呼ばれる場合もある）、又はデータ処理システムによって既に処理されているデータセット（中間データセットと呼ばれる場合もある）とすることができる。変換とは、データシンクに与えられるダウンストリームデータセットをもたすためにアップストリームデータセットに適用されるデータ処理操作である。データ系列図は、データ処理システム内のデータ系列要素のグラフィカル表現である。

【0030】

図1は、データ処理システムによって生成される出力データ110に関するデータ系列図100の一例である。図1の例では、データ処理システムがソースデータ102、104の2つのセットを受信する。ソースデータは、例えば計算システム内のデータを記憶するための単層ファイル等のファイル、リレーショナルデータベース若しくはオブジェクトデータベース等のデータベース、待ち行列若しくは別のリポジトリ内に記憶されるデータ記録、又はそれらから受信されるデータ記録とすることができる。例えば、ソースデータ102は、ファイル「US_feed.dat」内に記憶された米国内でのクレジットカード取引のデータ記録とすることができる。各データ記録は、記録構造内で定められる属性又はデータベーステーブル内のカラム等、1つ又は複数のフィールドのそれぞれのための値を含むことができる。ソースデータ102、104はバッチ単位で受信され処理され得る（例えば、毎時、毎日、毎週、毎月、毎四半期、毎年、又は別の間隔で処理されるファイル又はデータベースからのデータ）。ソースデータ102、104はストリームとして受信され、継続的に処理されてもよく、例えば待ち行列によってバッファされ、データが入手可能であり且つシステム資源が許すときに処理され得る。

【0031】

ソースデータ102は変換要素106によって処理され、変換要素106は、例えばソースデータ102を何らかの方法で変えるためにソースデータ102に作用する。変換要素は、仮想マシン内で実行されるjavaプログラム、実行ファイル、データフローグラフ、別の種類の実行可能プログラム等、データを操作可能な実行可能プログラムとすることができる。例えば、変換要素106は「TransformA.exe」と名付けられた実行ファイルであり得る。ある具体例では、変換要素106が、正しくないフォーマットを有するデータ記録等、ソースデータ102から不所望のデータ記録をフィルタで除去するフィルタコンポーネントであり得る。変換要素106は、基準データ120を考慮してソースデータ10

10

20

30

40

50

2 を処理して中間データ 1 1 2 をもたらす。基準データは、変換要素がデータを処理することを可能にするために変換要素によって使用されるデータである。例えば、マッピング操作を可能にする基準データは、処理されているデータ内の 1 つ又は複数のフィールド内の値に対応する値を有する 1 つ又は複数のフィールドを含む。中間データ 1 1 2 は、計算システム内のデータを記憶するためのファイル、データベース、待ち行列、又は別のリポジトリ内に記憶され得る。

【 0 0 3 2 】

変換要素 1 0 8 は、基準データ 1 2 2 を考慮してソースデータ 1 0 4 のセットを処理して中間データ 1 1 4 をもたらす。中間データ 1 1 4 は、計算システム内のデータを記憶するためのファイル、データベース、待ち行列、又は別のリポジトリ内に記憶され得る。

10

【 0 0 3 3 】

中間データ 1 1 2、1 1 4 は、基準データ 1 1 8 を使用する変換要素 1 1 6 によって一緒に処理される。一例では、変換要素 1 1 6 がマッピング操作であり、基準データ 1 1 8 が州の値及び対応する地域の値を示すデータ記録を含む。中間データ 1 1 2、1 1 4 が変換要素 1 1 6 によって処理されるとき、中間データ 1 1 2、1 1 4 内の各データ記録内の州フィールド内の値が基準データ 1 1 8 内で示される対応する地域にマップされる。一例では、基準データ 1 1 8 が、法人企業体、対応する部門識別子、経営者名、及び位置を示すビジネスデータを含む。中間データ 1 1 2、1 1 4 が変換要素 1 1 6 によって処理されるとき、基準データセットによって可能にされるマッピングに基づいて各データ記録が法人企業体に割り振られる。基準データ 1 1 8 は複数のデータセットを処理するために使用

20

【 0 0 3 4 】

変換要素 1 1 6 は、計算システム内のデータを記憶するためのファイル、データベース、待ち行列、又は別のリポジトリ内に記憶される出力データ 1 1 0 を出力する。出力データ 1 1 0 は、例えば同じデータ処理システム内の若しくは異なるデータ処理システム内の他の変換要素によって更に処理されてもよく、又は将来分析するために記憶され得る。

【 0 0 3 5 】

図 1 の例では、単一のデータ処理システム内のデータ系列要素に関して出力データ 1 1 0 のデータ系列が図示されている。一部の例では、複数のデータ処理システムによってデータセットのデータ系列を追跡することができる。例えば、出力データ X をもたらしするためにソースデータが第 1 のデータ処理システムによって最初に処理され得る。第 2 のデータ処理システムが、第 1 のデータ処理システムからの出力データ X を読み取り、出力データ X を処理して出力データ Y を生成する。出力データ Y は第 3 のデータ処理システムによって処理され、出力データ Z が生成される。出力データ Z のデータ系列は、最初のソースデータ、3 つのデータ処理システムのそれぞれに含まれる変換、及び 3 つのデータ処理システムの何れかによる処理中に使用される任意の基準データを含む。

30

【 0 0 3 6 】

一部の例では、目標要素 2 0 6 A のための終端間データ系列図 2 0 0 A の例で図示されているような、より複雑なデータ処理システムによって出力データが生成され得る。データ系列図 2 0 0 A では、データ要素 2 0 2 A と変換要素 2 0 4 A との間のつながりが図示されている。データ要素 2 0 2 A はデータセット、データセット内のテーブル、テーブル内のカラム、ファイル内のフィールド、又は他のデータを表し得る。変換要素の一例は、データ要素の単一出力がどのように作り出されるかを記述する実行ファイルの要素である。図 2 のデータ処理システム内で目標要素 2 0 6 A 内の（又は別のデータ要素 2 0 2 A 内の）潜在的なデータ品質問題の根本的原因を追跡することができる。図 2 の更なる説明は、参照によりその全内容を本明細書に援用する米国特許出願公開第 2 0 1 0 / 0 1 3 8 4 3 1 号に見出すことができる。

40

【 0 0 3 7 】

図 1 又は図 2 のデータ系列図等のデータ系列図内で示されている情報は、何れのアップ

50

ストリームデータソース、データシンク、又は変換がダウンストリームデータに影響を及ぼすかを示す。例えば、図1のデータ系列図100は、出力データ110がソースデータ102、104、基準データ118、120、122、及び変換要素106、108、116の影響を受けていることを明らかにする。

【0038】

ダウンストリームデータセット（出力データ110等）の系列を理解することは、ダウンストリームデータ内で生じ得るデータ品質問題の根本的原因を識別する際に有用であり得る。データ品質問題の根本的原因とは、ダウンストリームデータ内のデータ品質問題の少なくとも部分的な原因であるアップストリームのシステム、操作、又はデータセットの識別を意味する。出力データ110等におけるダウンストリームデータセット内のデータ品質問題の原因は、低品質のソースデータ、低品質の基準データ、出力データ110のセットのアップストリーム系列内の変換要素内の誤り、又はそれらのうちの何れか2つ以上の組合せであり得る。データ系列要素の品質又は状態を追跡することは、低品質の出力データのあり得る根本的原因を評価するために使用できる情報をもたらし得る。

10

【0039】

データセットのデータ品質とは、そのデータセットが予期される特性を有するかどうかを概して意味する。低いデータ品質は、予期された通りに振る舞わない、例えば統計的標準を外れる、標準的な照会に応答してルックアップの失敗又は別の種類の挙動を返すデータセット内に現れ得る。データセットの品質は、以下で説明するように、データセット内のデータ記録の一部若しくは全てのプロファイルに基づいて、特定のデータ記録の1つ又は複数のフィールドのそれぞれの中の値に基づいて、又はその両方に基づいて特徴付けることができる。

20

【0040】

ダウンストリームデータセット（例えば、出力データ110）内の低いデータ品質は、出力データのアップストリームデータ系列内の様々な要因の何れかにさかのぼることができる。低品質の出力データの1つのあり得る原因は、低品質のソースデータ、低品質の基準データ、又はその両方であり得る。例えば、あるソースデータセットは、伝送中に破損し若しくは中断されている場合があり、間違っただデータセットである可能性があり、欠落データを有することがあり、又は別の問題を有し得る。基準データセットは、基準データセットに対する最近の更新で誤りにさらされている場合があり、破損していることがあり、間違っただデータセットである可能性があり、又は別の問題を有し得る。低品質の出力データの別のあり得る原因は、出力データのアップストリームデータ系列内の変換要素の問題であり得る。例えば、変換要素を実装するソフトウェアが新バージョンに最近更新された場合、例えば更新されたソフトウェアが誤りを有し又は破損している場合、変換要素はもはや所望の処理を行わない場合がある。出力データセット内で生じ得る潜在的なデータ品質問題を先制して識別すること、出力データセット内で生じたデータ品質問題の根本的原因を後に追跡すること、又はその両方を容易にするために、出力データ110のセットのデータ系列内のソースデータ、基準データ、及び変換要素をモニタすることができる。

30

【0041】

ソースデータ及び基準データをモニタリングし分析することは、低品質の出力データの1つ又は複数のあり得る原因をユーザが診断するのを補助し得る。例えば、低品質の出力データセットが生成される場合、所与のソースデータセット又は基準データセット自体が低品質かどうか、従って低品質の出力データのあり得る一因かどうかを、低品質の出力データセットのデータ系列内のソースデータ又は基準データを分析することが示し得る。ソースデータ及び基準データをモニタリングすることは、処理された場合にダウンストリーム出力データ内でデータ品質問題を引き起こし得る、低品質のソースデータ又は基準データを先制して識別することもできる。

40

【0042】

図3A及び図3Bは、図1に示したデータ系列を有する出力データ110のセット内の未知の又は潜在的なデータ品質問題の根本的原因を追跡する手法を示す。図3Aを参照す

50

ると、入力データ（例えば、図1のソースデータ102、104）を処理する前に、基準データ118、120、122の品質が品質要素154、156、158によってそれぞれ特徴付けられる。一部の例では、基準データセットが更新されるとき、予定時刻に（例えば、周期的に又は基準データの更新が予定されているとき）、各入力データセットを処理する前に、又は他の時点において基準データの品質を特徴付けることができる。

【0043】

データセットの品質を特徴付けるために、品質要素がデータセット内のフィールドのプロファイル（統計調査と呼ばれる場合もある）を計算する。データ記録セットのプロファイルは、データ記録内のデータ値の例えばフィールドごとの要約である。プロファイルは、セット内のデータ記録の少なくとも一部のそれぞれのうちの1つ又は複数のフィールドのそれぞれのうちのデータ値を特徴付ける統計、値のヒストグラム、最大値、最小値、平均（例えば、中間又は中央）値、平均値からの標準偏差、個別値の数、（例えば、データセットごとの重要なデータ要素に関する）1つ又は複数のフィールド内の最も高頻度の値及び最も低頻度の値の標本、又は他の統計を含む。一部の例では、プロファイルが、データ記録内の1つ又は複数のフィールドのそれぞれのうちのデータ値を特徴付ける処理された情報を含み得る。例えば、プロファイルは、フィールド内の値の分類（例えば、収入データフィールド内のデータの高い、中位、又は低いカテゴリへの分類）、個々のデータ記録内のデータフィールド間の関係の指示（例えば、州のデータフィールドとZIPのデータフィールドとが無関係ではないという指示）、データ記録間関係（例えば、顧客__識別子フィールド内で共通値を有するデータ記録は関係しているという指示）、又はデータ記録セット内のデータを特徴付ける他の情報を含み得る。

10

20

【0044】

次いで、品質要素が1つ又は複数の規則を適用して、データセット内の任意の実際又は潜在的なデータ品質問題を識別する。以下で更に論じるように、規則はユーザによって指定されてもよく、プロファイルの許容可能な特徴又は禁止された特徴を示し得る。ある具体例では、基準データセットが米国の州の略記を列挙するフィールドを含む場合、規則の一例は、そのフィールド内の個別値の数が50を上回る場合、データ品質問題を識別すべきであると示すことができる。一部の例では、規則がデータセットの履歴プロファイル、例えば履歴平均値に基づき得る。データ品質問題がデータセット内で識別されない場合、規則を更新するために、例えば履歴平均値を更新するためにデータセットのプロファイルを使用することができる。実際の又は潜在的なデータ品質問題を有するものとして基準データセットが識別される場合、データ品質問題が対処されるまで処理を休止することができる。

30

【0045】

図3Bを参照すると、ソースデータ102、104の品質は品質要素150、152によってそれぞれ特徴付けられる。品質要素150、152は、データ処理システム内にデータが受信されるとき、それぞれのソースデータの予定された処理の前に、又は他の時点においてソースデータ102、104のそれぞれのデータ品質を特徴付けることができる。既知の又は潜在的なデータ品質問題を有するものとしてソースデータセットが識別される場合、例えばユーザに警告するために又は将来参照するためにデータ記憶域内に記憶するために、データ品質問題に関する情報を出力することができる。例えば、各品質要素150、152が対応するデータセットからデータを読み取ると、品質要素150、152はデータセットのプロファイルを計算する。

40

【0046】

ある具体例では、ソースデータ102のプロファイルを計算するために、品質要素150はソースデータ102内の取引__量フィールド内の値の全ての和を計算することができる。ソースデータ102のための規則が、取引__量フィールド内の値の全ての和を過去30ランにわたるその和の平均及び標準偏差と比較することができ、ソースデータ102の取引__量フィールド内の値の全ての和が和の平均値からの1標準偏差の範囲外である場合、データ品質問題を識別すべきであると示すことができる。

50

【 0 0 4 7 】

一部の例では、データセットの品質を特徴付けるために使用される規則が、データセット内のデータ記録のプロファイルの許容可能な特徴又は禁止された特徴を示し得る。プロファイルの特徴は値又は値域とすることができる。プロファイルの許容可能な特徴を示す規則は、プロファイルが許容可能な特徴を含む場合に満たされる。フィールドに関する許容可能な特徴の一例は、そのフィールドの許容可能な最大値及び最小値とすることができる。そのフィールドの平均値が許容可能な最大値と最小値との間に含まれる場合に規則が満たされる。プロファイルの禁止された特徴を示す規則は、禁止された特徴をプロファイルが含めない限り満たされる。フィールドに関する禁止された特徴の一例は、そのフィールドについて禁止されている値の一覧とすることができる。そのフィールドが禁止値の何れかを含む場合、その規則は満たされない。

10

【 0 0 4 8 】

プロファイルの特徴を示す規則は、特定のデータセットのフィールドのプロファイルとデータセットのフィールドに関する基準プロファイルとの間の許容偏差を示し得る。対応する規則によって示される許容偏差を上回るデータセットのプロファイルとデータセットに関する基準プロファイルとの間の偏差は、そのデータセット内のデータ品質問題の指示、従ってそのデータセットがダウンストリームデータセット内の既存の又は潜在的なデータ品質問題のあり得る根本的原因であるという指示であり得る。一部の例では、最大許容値及び最小許容値等の値域として許容偏差を指定することができる。一部の例では、平均値（例えば、過去のデータセット内の値の中間又は中央）とすることができる単一値からの標準偏差として許容偏差を指定することができる。

20

【 0 0 4 9 】

一部の例では、データセットの品質を特徴付けるために使用される規則は、フィールド内の値の妥当性等に基づき、データ記録の1つ又は複数のフィールドのそれぞれの中の値の許容特性又は禁止特性を示すことができる。フィールドに関する許容特性を示す規則は、フィールド内の値がその許容特性に適合する場合に満たされる。フィールドに関する禁止特性を示す規則は、フィールド内の値が禁止特性に適合しない限り満たされる。規則を満たす値は有効値と呼ばれる場合もあり、規則を満たさない値は無効値と呼ばれる場合もある。規則による許容特性又は禁止特性としてフィールド内の値の様々な特性を示すことができる。規則の一例は、許容値域若しくは禁止値域、最大許容値、最小許容値、又は許容若しくは禁止されている1つ又は複数の特定の値の一覧等、フィールドのコンテンツの許容特性又は禁止特性を示すことができる。例えば、1900未満の値又は2016を上回る値を有する誕生__年フィールドは無効と見なされ得る。規則の一例は、フィールドのデータタイプの許容特性又は禁止特性を示すことができる。規則の一例は、特定のフィールド内に値がないこと（又はヌルがあること）が許容されているか、又は禁止されているかを示すことができる。例えば、文字列値（例えば、「Smith」）を含むラスト__ネームフィールドは有効と見なされ得る一方、空白であり又は数値を含むラスト__ネームフィールドは無効と見なされ得る。規則の一例は、同じデータ記録内の2つ以上のフィールド間の許容又は禁止された関係を示し得る。例えば、ある規則は、州フィールドのあり得る各値に対応するZIPフィールドの値の一覧を指定することができる。その一覧によってサポートされないZIPフィールドの値と、州フィールドの値とのいかなる組合せも無効であると指定することができる。

30

40

【 0 0 5 0 】

一部の例では、履歴データを自動で分析することに基づいて規則を生成することができる。この種の規則を自動生成規則と呼ぶ。自動生成規則は、データセット内のデータ記録のプロファイルの許容可能な特徴又は禁止された特徴を示すことができる。例えば、プロファイルの自動生成規則は、特定のデータセットのフィールドのプロファイルとデータセットのフィールドの自動的に決定された履歴基準プロファイルとの間の許容偏差を示し得る。データセットの履歴基準プロファイルは履歴データに基づくことができ、例えば、履歴基準プロファイルは前日の同じデータセットのプロファイル、過去の複数日の（例えば

50

、先週又は先月にわたる)同じデータセットの平均プロファイル、同じデータセットの継続期間の平均プロファイルとすることができる。より広義には、基準プロファイルは、様々な種類の統計的分析を活用するための多岐にわたる基準情報を保持することができる。例えば、基準プロファイルは、値分布の標準偏差又は他の指示に関する情報を含み得る。以下の例では、及び本願の一般的な概念を限定することなく、基準プロファイルが過去のデータセットの数値的平均及び場合により標準偏差も含むと仮定する。

【0051】

自動生成規則は、データ記録のフィールド内の値の自動的に決定された許容特性又は禁止特性を示すことができる。一例では、あるフィールドに関する自動生成規則が、そのフィールドの履歴的な最大値又は最小値の分析に基づくフィールドの許容可能な最大値又は最小値を示し得る。一例では、あるフィールドに関する自動生成規則が、そのフィールドについて過去に生じた値の分析に基づくフィールドの許容値一覧を示し得る。一部の例では、データセットの全フィールドについて自動生成規則が指定される。一部の例では、フィールドの部分集合について規則が指定される。規則が指定されるフィールドは、例えばデータ記録の分析に基づいて自動で識別され得る。例えば、自動生成規則を生成可能なフィールドとして、少数の個別値を概して有するデータ記録セット内の任意のフィールド(低濃度フィールドと呼ばれる場合もある)を識別することができる。

10

【0052】

一部の例では、自動生成規則を生成するために機械学習技法が使用される。例えば、規則を生成する前に履歴平均又は期待値を識別するために、学習期間にわたってデータを分析することができる。学習期間は、指定の期間とすることができ、又は平均値若しくは期待値が安定した値に収束するまでの時間とすることができ。

20

【0053】

一部の例では、規則がユーザによって指定され得る。この種の規則をユーザ指定規則と呼ぶ。ユーザ指定規則は、特定のデータセットのフィールドのプロファイルの許容特性又は禁止特性、データセット内のデータ記録の1つ又は複数のフィールドのそれぞれの中の値の許容特性又は禁止特性、又はその両方を指定することができる。ユーザは、例えばシステムによって処理されるデータ記録の予期される特性についての自らの理解に基づいて規則を指定することができる。一部の例では、ユーザによって修正され得る省略時値をユーザ指定規則に割り当てることができる。

30

【0054】

ある具体例では、ソースデータが米国内で生じる取引に関するクレジットカード取引記録である。ソースデータは、1時間のインクリメント単位で処理されるストリーミングデータである。ソースデータについての、及びクレジットカード取引記録を処理するときに行うべき操作についての自らの知識に基づき、ユーザは、プロファイルすべき重要なデータ要素として取引識別子フィールド、カード識別子フィールド、州フィールド、日付フィールド、及び金額フィールドを識別することができる。

【0055】

ソースデータがクレジットカード取引記録である具体例では、ユーザは、州フィールドについて50個の許容値のみがあることを知っている場合がある。ユーザは、基準に対するソースデータセットのプロファイルの標準偏差に関係なく、ソースデータセットのプロファイルが州フィールド内で50を上回る値を識別する場合、警告フラグの使用を引き起こす規則を作成することができる。ユーザは、処理と同日に完了した取引に関するクレジットカード取引記録のみがソースデータセット内にあることも知っている場合がある。ユーザは、任意のソースデータ記録が処理日と一致しない日付を有する場合、警告メッセージの送信を引き起こす規則を作成することができる。

40

【0056】

図4を参照すると、一部の例では、ユーザが1つ又は複数の規則をユーザインタフェース400によって指定することができる。ユーザインタフェース400の一例は、複数の行402及び複数の列404を含む。各行402は、データセット内のデータ記録のフィ

50

ールド406に関連し、各列404は規則408に関連する。ユーザインタフェース400により、ユーザは1つ又は複数のフィールド406に関する規則を指定することができ、又はフィールドに関する事前にデータ投入された省略時規則を承認することができる。ユーザインタフェース400についての更なる説明は、参照によりその全内容を本明細書に援用する、2012年10月17日に出願された米国特許出願第13/653,995号に見出すことができる。ユーザインタフェース400の他の実装形態もあり得る。

【0057】

一部の例では、基準データセットの新バージョン内又はソースデータセット内等、起こり得るデータ品質問題がデータセット内で検出される場合、データベース内に記憶される根本的原因データセットの一覧上に起こり得るデータ品質問題を有するデータセットの識別子が配置される。出力データ110のセットのデータ品質問題が後に検出される場合、出力データ110のセットのアップストリームデータ系列要素を識別し、(存在する場合には)それらのアップストリームデータ系列要素の何れが根本的原因データセットの一覧に含まれているかを判定するためにデータベースを照会することができる。

10

【0058】

一部の例では、基準データセットの新バージョン又はソースデータセット内等、起こり得るデータ品質問題がデータセット内で検出される場合、ユーザ通知を使用可能にすることができる。一部の例では、データ品質問題を示すために警告フラグを記憶することができる。例えば、起こり得るデータ品質問題が基準データセットの新バージョン内で検出される場合、基準データの新バージョンに関するプロファイルデータと共に警告フラグを記憶することができる。起こり得るデータ品質問題がソースデータセット内で検出される場合、そのソースデータセットに関するプロファイルデータと共に警告フラグを記憶することができる。一部の例では、起こり得るデータ品質問題の存在を示すために警告メッセージをユーザに伝達することができる。警告メッセージは、例えばユーザインタフェース上のメッセージ、アイコン、ポップアップウィンドウ、電子メール、ショートメッセージサービス(SMS)メッセージ、又は別の形式とすることができる。

20

【0059】

一部の例では、警告フラグ又は警告メッセージが使用される基準プロファイルからの1つ又は複数の限界偏差を規則が指定することができる。例えば、現在のデータセットのプロファイルとそのデータセットに関する基準プロファイルとの間の偏差が、1~2標準偏差等だけ小さい場合に警告フラグを記憶することができ、偏差が2標準偏差を上回る場合に警告メッセージを伝達することができる。限界偏差は、各ソースデータセット及び基準データセットに固有であり得る。

30

【0060】

偏差が極度である、例えば基準プロファイルを3標準偏差だけ上回る場合等の一部の例では、ユーザが介入するまでデータ処理システムによる更なる処理を停止することができる。例えば、極度の偏差を有するソースデータ又は基準データの影響を受ける任意の更なる処理が一時停止される。一時停止すべき変換は、影響を受けるソースデータ又は基準データのダウストリームにあるデータ系列要素を参照するデータによって識別され得る。

40

【0061】

一部の例では、基準プロファイルデータが自動で決定される。例えば、所与のデータセットの基準プロファイルデータをそのデータセットの過去のプロファイルデータの履歴的な移動平均として(例えば、そのデータセットの新たなプロファイルデータが決定されるたびに基準プロファイルデータを再計算することによって)自動で更新することができる。一部の例では、ユーザが、例えば、所望の特性を有するデータセットをプロファイリングすることにより、最初の基準プロファイルデータを供給することができる。

【0062】

変換要素106、108、116のそれぞれに対する最近の更新の時間又は日付等、出力データのデータ系列内の変換要素106、108、116の更新状態を追跡することができる。変換要素に対する最新の更新のタイミングにアクセスすることができ、ユーザは

50

変換要素の1つ又は複数、例えば正しくない又は破損した変換要素が、出力データ110内の既存の又は潜在的なデータ品質問題のあり得る根本的原因かどうかを評価することができる。例えば、出力データ110が変換要素116から出力される直前に変換要素116が更新された場合、出力データ110内の既存の又は潜在的なデータ品質問題のあり得る根本的原因として変換要素116を識別することができる。

【0063】

図5を参照すると、追跡エンジン500は、ソースデータ及び基準データ等のデータ系列要素のプロファイル、並びにデータ処理システムによって生成される出力データ等の所与のデータセットのアップストリームデータ系列内の基準データ及び変換等のデータ系列要素に対する更新をモニタする。

【0064】

追跡エンジン500は、データ処理システムによって生成される出力データ等の所与のデータセットのアップストリームにあるデータ系列要素を参照するデータ504を記憶する、データ系列リポジトリ502を含む。例えば、データ系列リポジトリ502は、各データ系列要素の識別子及びデータ系列要素間の関係を示すデータを記憶し得る。データ系列リポジトリ502は、ファイル、データベース、又は別のデータ記憶機構であり得る。

【0065】

追跡エンジン500は、更新モニタ506を含む。更新モニタ506は、データ処理システム内の変換要素及び基準データセットが何れの時点で更新されるかをモニタする。データ系列リポジトリ502によって参照される変換要素ごとに、更新モニタ506は、変換要素を実装するソフトウェアが何れの時点で更新されるかをモニタする。更新が生じるとき、更新モニタ506は、ファイル、データベース、別のデータ記憶機構等の更新リポジトリ508内にエントリ510を記憶する。エントリ510は、ソフトウェアが更新された日付、時間、その両方等の更新のタイミングを示す。一部の例では、エントリ510が、更新についての手入力された説明、更新によって変更された命令行のテキスト、更新の性質についての別の指示等、更新の性質についての指示も含み得る。更新リポジトリ508は、変換要素の識別子により、更新のタイミングにより、又はその両方により索引を付けられ得る。

【0066】

データ系列リポジトリ502によって参照される基準データセットごとに、更新モニタ506は、基準データセットが何れの時点で更新されるかをモニタする。更新が生じるとき、更新モニタ506は、ファイル、データベース、別のデータ記憶機構等のプロファイルリポジトリ516内にエントリ514を記憶する。エントリ514は、基準データセットが更新された日付、時間、その両方等の更新のタイミングを示す。プロファイルリポジトリ516は、基準データセットの識別子により、更新のタイミングにより、又はその両方により索引を付けられ得る。

【0067】

基準データセットが更新されるとき、その基準データセットに関する品質要素が、基準データの新しいバージョンと呼ばれる場合もある更新された基準データのプロファイルを生成する。品質要素は、ファイル、データベース、別の記憶機構等の規則リポジトリ522内に記憶される重要なデータ要素の一覧520に従ってプロファイルを生成することができる。重要なデータ要素とは、ユーザ又はシステムにとって重要であることが分かっているデータ記録内のフィールド、例えばユーザによって指定されるフィールド又は自動で識別されるフィールドである。基準データの新しいバージョンに関する重要なデータ要素ごとにプロファイルが生成される。例えば、所与の重要なデータ要素について生成されるプロファイルは、重要なデータ要素に関する幾つの個別値が基準データセット内にあるか、及び各個別値が発生する回数を示す統計調査データであり得る。それぞれの重要なデータ要素の生成プロファイルを示す基準プロファイルデータ524が、例えば基準データに対する更新を示すエントリ514に関連してプロファイルリポジトリ516内に記憶される。

【0068】

10

20

30

40

50

ソースデータがデータ処理アプリケーションに与えられるとき、データ系列リポジトリ 5 0 2 によって参照される各ソースデータセットのプロファイルが対応する品質要素によって生成される。ソースデータ内の重要なデータ要素ごとにプロファイルが生成され、重要なデータ要素は、規則リポジトリ 5 2 2 内に記憶される重要なデータ要素の一覧 5 2 0 内で指定される。プロファイルされた各ソースデータセットの生成プロファイルを示すソースプロファイルデータ 5 2 6 が、ファイル、データベース、別のデータ記憶機構等のプロファイルリポジトリ 5 1 6 内に記憶される。

【 0 0 6 9 】

一部の例では、ダウンストリーム出力データ内でデータ品質問題が生じる場合にのみ、基準プロファイルデータ 5 2 4 及びソースプロファイルデータ 5 2 6 がアクセスされる。一部の例では、基準データの新バージョン又は受信されたソースデータそれぞれの潜在的なデータ品質問題をデータが示すかどうかを判定するために、基準プロファイルデータ 5 2 4、ソースプロファイルデータ 5 2 6、又はその両方がプロファイルモジュールによって分析される。プロファイルデータ 5 2 4、5 2 6 は、プロファイルの生成直後に分析することができ、又は後の時点において、例えば追跡エンジンが分析のために空いた計算資源を有する任意の時点において分析することができる。

10

【 0 0 7 0 】

基準プロファイルデータ 5 2 4 又はソースプロファイルデータ 5 2 6 を分析するために、分析モジュール 5 3 0 が、規則リポジトリ 5 2 2 内に記憶される自動生成規則又はユーザ指定規則等の規則 5 3 6 を適用する。規則は、例えばデータセットごとの 1 つ又は複数の重要なデータ要素、データ品質問題を引き起こし得る限界偏差、又は別の種類の規則を示し得る。

20

【 0 0 7 1 】

一部の例では、潜在的なデータ品質問題が基準データの新バージョン内又はソースデータセット内で検出される場合、データ系列リポジトリ 5 0 2 内に記憶される根本的原因データセットの一覧 5 5 0 上に潜在的なデータ品質問題を有するデータセットの識別子が配置される。ユーザがダウンストリームデータセットのデータ品質問題を後に検出する場合、出力データセットのアップストリームにあるデータ系列要素を識別し、(存在する場合には)それらのアップストリームデータ系列要素の何れが根本的原因データセットの一覧 5 5 0 上に含まれているかを識別するために、ユーザは、データ系列リポジトリ 5 0 2 を照会することができる。

30

【 0 0 7 2 】

一部の例では、起こり得るデータ品質問題があるかどうかを判定するために出力データ 1 1 0 が自動で分析される。例えば、現在の出力データ 1 1 0 のプロファイルを出力データ 1 1 0 の旧バージョンの基準プロファイルと比較するために、例えば出力データ 1 1 0 の各バッチ又は時間間隔をプロファイルすることができ、プロファイリング規則及び検証規則を出力データ 1 1 0 に適用することができる。出力データプロファイリング規則内で指定されるように、現在の出力データ 1 1 0 のプロファイルが基準プロファイルから閾値量を上回って外れる場合、潜在的なデータ品質問題を有するものとして現在の出力データ 1 1 0 を識別することができる。出力データ検証規則内で指定されるように、現在の出力データ 1 1 0 内の特定のデータ要素が期待値域から閾値量を上回って外れる値を有する場合、潜在的なデータ品質問題を有するものとして現在の出力データ 1 1 0 を識別することができる。データウェアハウス内に警告フラグを出力データ 1 1 0 と共に記憶することができ、又はユーザは、例えば、ユーザインタフェース又はメッセージによって通知され得る。

40

【 0 0 7 3 】

一部の例では、潜在的なデータ品質問題を有するものとしてユーザが出力データ 1 1 0 のセットを識別する。例えば、複数の出力データ 1 1 0 のセットを要約するレポートを作成するビジネスアナリストは、自らが分析している他の出力データセットに比べて特定の出力データ 1 1 0 のセットが殆ど意味をなさないことを認識する場合がある。アナリスト

50

は、潜在的なデータ品質問題を有するものとしてその特定の出力データ 110 のセットのフラグを立てることができる。

【0074】

出力データがデータ品質問題を有する場合、データ品質問題の根本的原因を識別するために、追跡エンジン 500 内に記憶されている情報がアクセスされ得る。例えば、ファイル名又はタイムスタンプ等の出力データの識別子が、例えば自動で又はユーザによって照会モジュール 548 に与えられ得る。照会モジュール 548 は、識別された出力データに関連し得る情報を求めて関連リポジトリのそれぞれを照会する。具体的には、照会モジュール 548 は、識別された出力データが依拠する変換、ソースデータ、及び基準データを識別するためにデータ系列リポジトリ 502 を照会する。次いで、照会モジュール 548 は、出力データを処理する直前に生じた識別された変換要素の何れかに対する更新を示す任意のエントリ 510 を求めて更新リポジトリを照会することができる。照会モジュール 548 は、関連する基準プロファイルデータ 524 及び関連する任意の警告フラグと共に、識別された基準データに対する更新を示す任意のエントリ 514 を求めてプロファイルリポジトリ 516 を照会することができる。照会モジュール 548 は、識別された任意のソースデータセットに関するソースプロファイルデータ 526 を求めてプロファイルリポジトリ 516 を照会することができる。

10

【0075】

照会モジュール 548 による照会に応じて返される結果がユーザインタフェース上で表示される。ディスプレイは、出力データ内のデータ品質問題の潜在的な根本的原因の理解を得るために、ユーザがデータを見て操作することを可能にする。例えば、出力データが処理される直前に変換要素に対するソフトウェア更新があった場合、更新についての説明又は変更された命令行等、ユーザはその更新に関連する情報を見ることができる。基準プロファイルデータ又はソースプロファイルデータに関連する警告フラグがあった場合、ユーザはプロファイルデータを見ることができる。

20

【0076】

一部の例では、照会モジュール 548 によって返される結果は、潜在的なデータ品質問題を有する出力データに関する処理を変換要素が行う直前に変換要素に対する更新が行われたことを示し得る。これを最近更新された変換要素と呼ぶ場合がある。直前とは、例えば、処理の 10 分以内、1 時間以内、1 日以内、別の時間内等、設定された時間内を意味する。更新モニタ 506 は、最近更新された変換要素の 1 つ又は複数が出力データ内のデータ品質問題の潜在的な根本的原因かどうかを示し得る、最近更新された変換要素に関する追加情報を得ることができる。例えば、更新モニタ 506 は、最近更新された変換要素に関連する任意の処理アーティファクトを識別することができる。処理アーティファクトがあることは、最近更新された変換要素の潜在的な問題を示し得る。最近更新された変換要素に対する更新を更新ログが反映することを確実にするために、更新モニタ 506 は最近更新された変換要素に関連する更新ログを点検することができる。更新ログと最近更新された変換要素に対する更新を示すデータ 510 との間の不一致は、変換要素の潜在的な問題を示し得る。最近更新された変換要素の更新中に取り込まれている可能性がある潜在的な誤りを識別するために、更新モニタ 506 はチェックサム又は他のシステムデータを点検することができる。

30

40

【0077】

一部の例では、最近更新された変換要素の潜在的な問題が検出される場合、ユーザ通知を使用可能にすることができる。一部の例では、潜在的な問題を示すために警告フラグを、例えば更新を示すデータ 510 と共に更新リポジトリ 508 内に記憶することができる。一部の例では、最近更新された変換要素の潜在的な問題の存在を示すために通信モジュール 546 によって警告メッセージをユーザに伝達することができる。例えば、警告メッセージは、ユーザインタフェース上のメッセージ、アイコン、ポップアップウィンドウ、電子メール、SMS メッセージ、又は別の形式とすることができる。一部の例では、データ系列及びデータ品質の分析を粗粒度のデータ系列と呼ぶ場合があるデータセットのレベ

50

ルとすることができる。粗粒度のデータ系列は、ダウンストリームデータセットのデータ系列を見る。ダウンストリームデータセットを生成するために使用されるアップストリームデータセット及びアップストリーム変換要素は、ダウンストリームデータセットのデータ系列内にあると見なされる。一部の例では、データ系列及びデータ品質の分析を細粒度のデータ系列と呼ぶ場合がある個々のフィールドのレベルとすることができる。細粒度のデータ系列は、ダウンストリームデータセット内の特定のフィールドのデータ系列を見る。ダウンストリームデータセット内の特定のフィールドを生成するために使用されるアップストリーム変換要素及びアップストリームデータセット内のフィールドは、ダウンストリームデータセットのデータ系列内にあると見なされる。データ品質の分析についてここで説明する手法は、粗粒度のデータ系列及び細粒度のデータ系列の両方に関連して適用され得る。

10

【0078】

プロファイリングに関する更なる情報は、参照によりその全内容を本明細書に援用する「Data Profiling」という名称の米国特許第8,868,580号に見出すことができる。典型的には、データ記録はデータフィールドの組に関連し、各フィールドは各記録に関する特定の値(場合によりヌル値を含む)を有する。一部の例では、データセット内のデータ記録が決まった記録構造を有し、かかる記録構造内では各データ記録が同じフィールドを有する。一部の例では、データセット内のデータ記録が、例えば可変長ベクトル又は条件付きフィールドを含む可変記録構造を有する。一部の例では、プロファイルモジュール218が、データセット内のデータ記録に関する初期フォーマット情報をプロファイル要素150、152、154に与えることができる。初期フォーマット情報は、例えば個別値を表すビット数(例えば、16ビット)、記録フィールドに関連する値及びタグ又は区切り符号に関連する値を含む値の順序、ビットによって表わされる値の種類(例えば、文字列、符号付き/符号なし整数又は他の種類)、又は他のフォーマット情報を含み得る。フォーマット情報は、規則リポジトリ522内に記憶されるデータ操作言語(DML)ファイル内で指定され得る。プロファイル要素150、152、154は、SQLテーブル、XMLファイル、CSVファイル等の様々な共通データシステムフォーマットのデータを自動で解釈するために既定のDMLファイルを使用することができ、又は専用のデータシステムフォーマットを記述する、規則リポジトリ222から得られるDMLファイルを使用することができる。

20

30

【0079】

図6は、ユーザが出力データセット内の潜在的なデータ品質問題の根本的原因を調査することを可能にするユーザインタフェース300の一例を示す。ユーザインタフェース300により、ユーザは出力データセットの識別子302又は出力データ内の特定のデータ要素の識別子304を入力することができる。例えば、識別子302又は304は、潜在的なデータ品質問題を有する出力データセット又は特定のデータ要素を識別することができる。図6の例では、ユーザが出力データセット「Billing_records.dat.」を入力している。ユーザインタフェース300上に、出力データの識別されたセット又は識別されたデータ要素のアップストリームのデータ系列要素をグラフィカルに示す対話型データ系列図310が表示される。データ系列図310の例では、識別された出力データセットのアップストリームのデータ系列要素がソースデータ312、314の2つのセット、2つの変換要素316、318、及び1つの基準データ320のセットを含む。

40

【0080】

この例のソースデータ312、変換要素318、基準データ320等、起こり得るデータ品質問題を有するアップストリームデータ系列要素が警告フラグ324a、324b、324cでそれぞれ印付けされる。ユーザは、警告フラグ上でクリックし若しくはタップし、警告フラグ上にマウスポインタを重ね、又は他の方法で警告フラグを選択すること等によって警告フラグを選択し、関連する起こり得るデータ品質問題に関する情報にアクセスすることができる。データセットに関連する起こり得るデータ品質問題に関する情報は、プロファイルデータ、1つ又は複数のデータ要素に関する基準プロファイルデータ、プ

50

ロファイルデータの統計的分析の結果（基準プロフィールデータからのプロフィールデータの偏差等）、検証規則によって指定された許容値を満たさない値、又は他の情報を含み得る。変換要素に関連する起こり得るデータ品質問題に関する情報は、変換要素に対する直近の更新日、更新についての説明、更新に由来する符号の抜粋、又は他の情報を含み得る。一部の例では、警告フラグの1つをユーザが選択することに応じてデータ系列図上に情報バブルをオーバーレイすることができる。一部の例では、警告フラグの1つをユーザが選択することに応じて新たな画面を表示することができる。一部の例では、情報バブル又は新たな画面内に表示される情報に対話型とすることができ、それにより、ユーザは、情報片を選択することによって更なる詳細情報にアクセスすることができる。

【0081】

ユーザインタフェース300により、ユーザは規則エディタ328にアクセスすることもでき、規則エディタ328によってユーザはプロフィールリング規則、検証規則、又はその両方を追加し、削除し、又は修正することができる。例えば、ユーザは、データセットごとに重要なデータ要素を追加し、削除し、又は修正すること、潜在的なデータ品質問題の識別を引き起こす限界偏差を更新すること、プロフィールリング規則又は検証規則を新たなデータセットの受信時に自動で適用すべきかダウンロードのデータ品質問題を検出したときにのみ自動で適用すべきかを指定すること、又はプロフィールリング規則若しくは検証規則に対して他の変更を加えることができる。

【0082】

ある具体例では、データ処理システムが通話記録を処理して課金記録を生成する。各ソースデータ記録は通話を表し、通話の日付、時間、通話の持続時間、ダイヤル側の電話番号、着呼側の電話番号等のデータを記憶するフィールドを含む。請求書を生成するために、ソースデータ記録がバッチ処理内で毎月処理される。この例では、2015年5月の月に顧客アカウントの95%について請求書が生成されなかった。ユーザは、2015年5月の請求書を生成するために使用された出力データのアップストリームデータ系列内のデータ系列要素のプロファイル及びかかるデータ系列要素に対する更新に関する情報を要求した。基準ソースプロフィールデータがダイヤル側電話番号フィールド内の150万個～240万個の固有値の期待範囲を示した一方、ソースプロフィールデータは、2015年5月の請求書を生成するために使用されたソースデータ記録内のダイヤル側電話番号フィールドが10個の固有値のみを有したことを示した。ソースプロフィールデータのこの点検に基づき、ユーザはソースデータ記録が破損していると判定した。2015年5月の請求書を正しく生成するためにソースデータ記録が圧縮記憶域から取得され、再び処理された。

【0083】

別の具体例では、データ処理システムが企業内財務記録を処理し、各財務記録を企業の事業部に割り振る。企業の事業部に対する各財務記録の割り振りは、各記録内の部門識別子を企業基準データセットによって提供される企業の6つの事業部のうちの1つにマッピングすることによって行われる。企業基準データに関する基準プロフィールデータは、企業の事業部の数が過去10年にわたって一貫して6つであったことを示した。基準データは毎四半期更新される。直近の更新後、基準データがプロフィールされ、基準データ内の企業の事業部の数が60まで増えたことを示した。6つの事業部の基準からの更新された基準データのプロファイルの偏差は、システム管理者への警告メッセージの送信を引き起こすのに十分大きかった。加えて、基準データを調査し、必要に応じて訂正できるまでデータ処理システムによる更なる処理が一時停止された。

【0084】

図7を参照すると、ある具体例では、データ処理システム50が、2016年4月1日にthebostonshop.comで行われたオンライン購入の記録を含む入力データ58を処理する複数の変換要素52、54、56を含む。入力データ58の各記録は、州フィールドを含む複数のフィールドを含む。この例では、コンポーネント56は、入力データの州フィールド内の値に基づいて8つのファイル60a～60hのうちの1つに各データ記録を送信

10

20

30

40

50

する分割コンポーネントである。例えば、州フィールド内にM Aの値を有する記録はファイル6 0 aに送信され、値T Xを有する記録はファイル6 0 bに送信され、値C Aを有する記録はファイル6 0 cに送信され、値D Eを有する記録はファイル6 0 dに送信され、値N Yを有する記録はファイル6 0 eに送信され、値I Lを有する記録はファイル6 0 fに送信され、値R Iを有する記録はファイル6 0 gに送信され、他の任意の値を有する記録はファイル6 0 hに送信される。各ファイルに送信される記録の数が図7に示されている。図7の例では、各ファイルに送信される記録の数が期待範囲内にあり、従ってデータ品質の警告は生成されない。これは入力データ5 8が期待範囲に含まれるためである。

【0085】

入力データ5 8の品質が品質要素6 2によって特徴付けられる。品質要素6 2は入力データ5 8の州フィールドのプロファイルを生成し、州フィールドのプロファイルと入力データの州フィールドの基準プロファイルとの間の許容偏差を示す自動生成規則を適用する。基準プロファイルは、過去1年にわたってデータ処理システム5 0によって処理されたデータの平均プロファイルを表し、それを超えると潜在的なデータ品質問題が識別される許容偏差を示す。この例では、入力データ5 8のプロファイル内の州フィールド内の値分布が基準プロファイル内の値分布と10%を超えて異なる場合、潜在的なデータ品質問題を有するものとして入力データ5 8を識別すべきであると自動生成規則が示し、州フィールドの基準プロファイルは、10%の許容偏差と共に州フィールド内の以下の値分布：

M A : 6 %

T X : 25 %

C A : 33 %

D E : 3 %

N Y : 17 %

I L : 11 %

R I : 4 %

他の任意の値 : 1 %

を示し、図7を見れば分かるように、州フィールドの実際のプロファイルが基準プロファイルの10%の許容偏差に含まれており、従って入力データのデータ品質問題はない。

【0086】

図8 Aを参照すると、データ処理システム5 0の異常動作の一例では、入力データ5 5が、2016年4月2日にbostonshop.comで行われたオンライン購入の記録を含む。この例では、ファイル6 0 gに記録が送信されていない。データ処理システム5 0のオペレータはファイル6 0 gが空であることに気付く場合があり、又は空のファイルがダウンロードのデータ処理システムによる更なる処理内で誤りを引き起こす可能性がある。データ処理システム5 0のオペレータは、ファイル6 0 a ~ 6 0 hのデータ系列内のアップストリームデータ要素の品質を調査することにより、ファイル6 0 gに記録が送信されなかった根本的原因を追跡することができる。具体的には、入力データ5 5はファイル6 0 a ~ 6 0 hのアップストリームデータ系列に属する。

【0087】

図8 Bも参照すると、品質要素6 2は入力データ5 5の州フィールドの以下の実際のプロファイルを生成する：

M A : 6 %

T X : 25 . 1 %

C A : 32 . 7 %

D E : 2 . 9 %

N Y : 17 . 1 %

I L : 11 . 1 %

R I : 0 %

他の任意の値 : 5 . 1 %。

【0088】

10

20

30

40

50

入力データ 55 の州フィールドのプロファイルと州フィールドの基準プロファイルとの間の偏差により、入力データ 55 は潜在的なデータ品質問題を有するものとして識別され、潜在的なデータ品質問題を示すために警告フラグが記憶される。オペレータが空のファイル 60 g の根本的原因を追跡するとき、オペレータは潜在的なデータ品質問題が入力データ 55 内に存在したことを容易に認めることができる。例えば、入力データ 55 が破損していたかどうか、アップストリームのデータ処理システム内での入力データ 55 のそれまでの処理が偏差を引き起こしたかどうか、又は別の原因を判定するために、オペレータはその知識を使用して偏差の原因を調査することができる。例えば、図 8 C も参照すると、この例では実際の入力データ 55 の一部を見ることにより、オペレータは、値「RI」内の文字が「IR」と読めるように逆になっており、それらの記録がファイル 60 g 内ではなくファイル 60 h 内に記憶されることを引き起こしていると認識し得る。

【0089】

図 9 A を参照すると、データ処理システム 50 の異常動作の別の例では、入力データ 64 が、2016 年 4 月 3 日に thebostonshop.com で行われたオンライン購入の記録を含む。この例では、記録がファイル 60 a のみに送信されており、他のファイル 60 b ~ 60 h の何れにも送信されていない。データ処理システム 50 のオペレータは、ファイル 60 b ~ 60 h が空であることに気付く場合があり、又は空のファイルがダウンストリームのデータ処理システムによる更なる処理内で誤りを引き起こす可能性がある。

【0090】

図 9 B も参照すると、データ処理システムのオペレータは、ファイル 60 a ~ 60 h のデータ系列内のアップストリームデータ要素の品質を調査することにより、記録の全てがファイル 60 a に送信された根本的原因を追跡することができる。この例では、品質要素 62 が入力データ 64 の州フィールドの以下のプロファイルを生成する：

MA : 6 . 1 %

TX : 25 . 2 %

CA : 32 . 6 %

DE : 2 . 9 %

NY : 17 . 0 %

IL : 11 . 1 %

RI : 4 . 1 %

他の任意の値 : 1 %。

【0091】

入力データ 64 の州フィールドのプロファイルは州フィールドの基準プロファイルと一貫性があり、従って潜在的なデータ品質問題は識別されない。次いで、オペレータは、ファイル 60 a ~ 60 h のデータ系列内にある変換要素 52、54、56 の更新状態を調査することができる。例えば、オペレータは、入力データ 64 を処理する直前に変換要素 56 が更新されたと判定することができ、従って変換要素 56 が空のファイル 60 b ~ 60 h の根本的原因であり得る。

【0092】

図 10 A を参照すると、ある具体例では、データ処理システム 80 が、特定の塔によって扱われる携帯電話の通話に関する通話記録を含む入力データ 86 のストリームを処理する複数の変換要素 82、84 を含む。入力データ 86 の各記録は、電話__番号フィールドを含む複数のフィールドを含む。入力データ 86 は変換要素 82 によってフォーマットされ、その後、電話__番号フィールド内の値により変換要素 84 によってソートされ、待ち行列 88 内に出力され、そこから更なる処理のために第 2 のデータ処理システム 90 内に供給される。この例では、待ち行列 88 から第 2 のデータ処理システム 90 内に供給される記録の 25 % が処理誤りを引き起こす。データ処理システム 80 のオペレータは、待ち行列 88 のデータ系列内のアップストリームデータ要素の品質を調査することにより、これらの処理誤りの根本的原因を追跡することができる。

【0093】

10

20

30

40

50

入力データ 86 の品質が品質要素 90 によって特徴付けられ、フォーマット変換要素 82 から出力されるデータ 94 の品質が品質要素 92 によって特徴付けられる。品質要素 90、92 の両方は、電話__番号フィールド内の値が 10 桁の整数であるべきこと、及び記録の 3% 超が規則を満たさない場合、潜在的なデータ品質問題を識別すべきことを指定するユーザ生成規則を適用する。この例では、品質要素 90 は、データ 86 内の記録の 0.1% が電話__番号フィールド内で 11 桁の整数を有すると判定する。この記録のパーセンテージは 3% の閾値を下回るため、品質要素 90 は入力データ 86 のいかなる潜在的なデータ品質問題も識別しない。品質要素 92 は、電話__番号フィールド内に英数字値を有するものとしてデータ 94 内の記録の 25% を特徴付ける。データ 94 の一部の一例を図 10B に示す。データ 94 の潜在的なデータ品質問題を示すために警告フラグが記憶される。オペレータが処理誤りの根本的原因を追跡するとき、オペレータは、入力データ 86 内のデータ品質問題が識別されていないが、潜在的なデータ品質問題がデータ 94 内にあることを容易に認めることができる。

10

20

30

40

50

【0094】

図 11 を参照すると、ソースデータセットの品質を決定するためのプロセスの一例では、ソースデータセットがデータ処理アプリケーション内に受信される(400)。ソースデータセットのプロファイルが生成され記憶される(402)。ソースデータセットに関する 1 つ又は複数の規則が取得される(404)。ソースデータ又はソースデータのプロファイルが 1 つ又は複数の規則に従って分析される(406)。1 つ又は複数の規則がソースデータセットによって満たされない場合(408)、潜在的なデータ品質問題を示す警告がプロファイルデータと共に記憶され、ユーザに伝達され、又はその両方が行われ(410)、起こり得るデータ品質問題を有するデータセットの一覧にソースデータが追加される。1 つ又は複数の規則がソースデータによって満たされる場合(408)、ソースデータがデータ処理アプリケーションによって処理される(412)。規則によって指定される閾値又は許容値からの極度の偏差等に関する一部の事例では、ユーザによる介入が処理の再開を可能にするまで処理が一時停止される。処理中又は処理後、ユーザは、例えば、ダウンストリームのデータ品質問題の潜在的な根本的原因を調査するために、記憶されたプロファイルデータにアクセスする権限を与えられる。

【0095】

図 12 を参照すると、データ処理システム内の基準データの品質をモニタリングするためのプロセスの一例では、基準データセットがモニタされる(500)。基準データセットが更新されると、基準データの新バージョンのプロファイルが生成され記憶される(502)。例えば、プロファイルの生成は、基準データに対する予定された各更新後に行われ得る。基準データセットに関する 1 つ又は複数の規則が取得される(504)。基準データの新バージョン又は基準データの新バージョンのプロファイルが 1 つ又は複数の規則に従って分析される(506)。1 つ又は複数の規則が基準データの新バージョンによって満たされない場合(508)、起こり得るデータ品質問題を示す警告がプロファイルデータと共に記憶され、ユーザに伝達され、又はその両方が行われる(510)。1 つ又は複数の規則が基準データの新バージョンによって満たされる場合(508)、データ処理システムによるその後の処理の開始又は続行が許可される(512)。規則によって指定される閾値又は許容値からの極度の偏差等に関する一部の事例では、ユーザによる介入が処理の開始又は続行を可能にするまで処理が一時停止される。処理中又は処理後、ユーザは、例えば、ダウンストリームのデータ品質問題の潜在的な根本的原因を調査するために、記憶されたプロファイルデータにアクセスする権限を与えられる。

【0096】

一部の例では、例えば規則ごとに更新日を決定するために、規則を適用する前に規則が分析される。規則が限界有効期間よりも古い場合、その規則は適用されない場合があり、又はその規則を更新する準備が整っている可能性があることをユーザに警告することができる。

【0097】

図13を参照すると、変換要素に対する更新を分析するためのプロセスの一例では、変換要素に対する最近の更新の時間が識別される(600)。例えば、最近の更新のタイムスタンプをデータリポジトリ内に記憶することができる。変換要素が最近の更新を有さない場合(602)、変換要素に対する更新が更に分析されることはない(604)。最近の更新は、10分以内、1時間以内、1日以内、別の時間内等、ある限界時間内の更新とすることができる。変換要素が最近更新された場合(602)、任意の処理アーティファクトが識別される(606)。更新ログとデータリポジトリ内に記憶された最近の更新のタイムスタンプとの間の任意の不整合を識別するために、変換要素に関連する更新ログが点検される(608)。変換要素の更新中に取り込まれている可能性がある任意の潜在的な誤りの指示のために、変換要素に関連するチェックサム又は他のシステムデータが点検される(610)。潜在的な問題が識別されない場合(612)、システムによる処理の開始又は続行が許可される(614)。1つ又は複数の潜在的な問題が識別される場合(612)、変換要素の潜在的な問題を示す警告がデータリポジトリ内に記憶され、ユーザに伝達され、又はその両方が行われる(616)。データ処理システムによる処理は、開始すること若しくは続行することを認められてもよく、又はユーザによる介入が処理の開始又は続行を可能にするまで一時停止してもよい。

10

【0098】

図14は、プロセスの一例のフローチャートである。データ処理システムによって生成される出力データセットを示す情報が受信される(700)。出力データセットに関係するデータ系列情報に基づき、出力データセットが依拠する1つ又は複数のアップストリームデータセットが識別される(702)。データ系列情報は、出力データセットが依拠する1つ又は複数のデータセット、出力データセットに依拠する1つ又は複数のデータセット、又はその両方を示す。1つ又は複数のデータセットの何れが誤り又は起こり得る誤りを有するかを判定することを含め、データセットの部分集合を識別するために出力データセットが依拠する識別されたアップストリームデータセットのそれぞれを分析する(704)。特定のアップストリームデータセットごとに、特定のアップストリームデータセットのプロファイルと特定のアップストリームデータセットに関する基準プロファイルとの間の許容偏差を示す第1の規則が適用され(706)、特定のアップストリームデータセット内の1つ又は複数のデータ要素の許容値又は禁止値を示す第2の規則が適用される(708)。一部の例では、第1の規則又は第2の規則のみが適用される。第1の規則、第2の規則、又はその両方は自動で生成され又はユーザによって指定され得る。第1の規則、第2の規則、又はその両方を適用した結果に基づき、アップストリームデータセットの1つ又は複数部分が部分集合のために選択される(710)。アップストリームデータセットの部分集合に関連する情報が出力される(712)。

20

30

【0099】

図15は、プロセスの一例のフローチャートである。データ処理システムのダウンストリームデータセットのデータ要素内の誤り又は起こり得る誤りが、例えば自動で又はユーザ入力に基づいて識別される(900)。ダウンストリームデータセットに関係するデータ系列情報に基づき、データ要素に影響を及ぼす1つ又は複数のアップストリームデータセットが自動で識別される(902)。識別されたアップストリームデータセットのそれぞれの現在のプロファイル及び基準プロファイルを分析することを含め、何れのアップストリームデータセットが誤りを有するか又は誤りを有する可能性が高いかを判定する(904)。例えば、各アップストリームデータセットは、現在のプロファイルのそれぞれに1つ又は複数の規則を適用することによって分析され得る。規則は、特定のアップストリームデータセットの現在のプロファイルと特定のアップストリームデータセットの対応する基準プロファイルとの間の許容偏差を示し得る。規則は、特定のアップストリームデータセット内のデータ要素に関する許容値を示し得る。誤りを有するか又は誤りを有する可能性が高いアップストリームデータセットのそれぞれに関連する情報が出力される(906)。

40

【0100】

50

本明細書に記載のデータ品質をモニタし追跡するための技法は、コンピュータ技術に基づいており、コンピュータによって実装されるプロセスの実行中に生じる問題に対処するために使用することができる。例えば、本明細書に記載のモニタリング及び追跡のための技法を使用し、コンピュータによって実装されるデータ処理システムによるデータセットの処理がモニタされ、より効率的、効果的、又は正確にされ得る。加えて、本明細書に記載の技法は、システム管理者等のユーザがデータ処理システムの動作を管理することを補助するために適用され得る。

【0101】

図16は、モニタリング及び追跡のための技法が使用され得るデータ処理システム1000の一例を示す。システム1000はデータソース1002を含み、データソース1002は、そのそれぞれが様々なフォーマット（例えば、データベーステーブル、スプレッドシートファイル、フラットテキストファイル、又はメインフレームコンピュータによって使用されるネイティブフォーマット）の何れかによってデータを記憶し又は提供することができる、記憶装置又はオンラインデータストリームへの接続等の1つ又は複数のデータソースを含み得る。データは、ロジスティックデータ、分析データ、又はマシンデータとすることができる。実行環境1004は、前処理モジュール1006及び実行モジュール1012を含む。実行環境1004は、例えばUNIXオペレーティングシステムのバージョン等の適切なオペレーティングシステムの制御下で1つ又は複数の汎用コンピュータ上にホストされ得る。例えば、実行環境1004は、ローカルである（例えば、対称型マルチプロセッシング（SMP）コンピュータ等のマルチプロセッサシステム）若しくはローカル分散された（例えば、クラスタとして結合される複数のプロセッサ若しくは超並列処理（MPP）システム、又は遠隔である若しくは遠隔分散された（例えば、ローカルエリアネットワーク（LAN）及び/又は広域ネットワーク（WAN）によって結合される複数のプロセッサ）、又はその任意の組合せである、複数の中央処理装置（CPU）又はプロセッサコアを使用するコンピュータシステムの構成を含むマルチノード並列計算環境を含み得る。

【0102】

データソース1002を提供する記憶装置は、例えば実行環境1004をホストするコンピュータに接続される記憶媒体（例えば、ハードドライブ1008）上に記憶され実行環境1004にとってローカルである場合があり、又は例えば実行環境1004をホストするコンピュータと（例えばクラウドコンピューティングインフラによって提供される）遠隔接続上で通信する遠隔システム（例えば、メインフレームコンピュータ1010）上にホストされ実行環境1004にとって遠隔的である場合もある。

【0103】

前処理モジュール1006は、データソース1002からデータを読み取り、実行用のデータ処理アプリケーションを作成する。例えば、前処理モジュール1006はデータ処理アプリケーションをコンパイルし、実行環境1004にとってアクセス可能なデータ記憶システム1016との間でコンパイルされたデータ処理アプリケーションを記憶及び/又はロードし、実行用のデータ処理アプリケーションを作成するための他のタスクを実行することができる。

【0104】

実行モジュール1012は、前処理モジュール1006によって作成されたデータ処理アプリケーションを実行してデータセットを処理し、その処理から生じる出力データ1014を生成する。出力データ1014はデータソース1002内に再び記憶することができ、実行環境1004にとってアクセス可能なデータ記憶システム1016内に記憶することができ、又は他の方法で使用され得る。データ記憶システム1016は、実行モジュール1012によって実行されるデータ処理アプリケーションを開発者1020が設計し編集することができる開発環境1018にとってもアクセス可能である。一部の実装形態では、開発環境1018は頂点間の（作業要素、即ちデータのフローを表す）有向辺によってつながれる（データ処理コンポーネント又はデータセットを表す）頂点を含むデータ

10

20

30

40

50

フローグラフとしてアプリケーションを開発するためのシステムである。例えば、かかる環境は、参照により本明細書に援用する「Managing Parameters for Graph-Based Applications」という名称の米国特許出願公開第2007/0011668号でより詳細に説明されている。かかるグラフベースの計算を実行するためのシステムは、参照によりその全内容を本明細書に援用する「EXECUTING COMPUTATIONS EXPRESSED AS GRAPHS」という名称の米国特許第5,966,072号で説明されている。このシステムに従って作成されるデータフローグラフは、グラフコンポーネントによって表わされる個々のプロセスの内外に情報を出し入れするための、プロセス間で情報を移動させるための、及びプロセスの実行順序を定めるための方法を提供する。このシステムは、使用可能な任意の方法からプロセス間通信方法を選択するアルゴリズムを含む（例えば、グラフのリンクによる通信経路は、プロセス間でデータを渡すためにTCP/IP又はUNIXドメインソケットを使用することができ、又は共用メモリを使用する）。

10

【0105】

前処理モジュール1006は、様々な形式のデータベースシステムを含むデータソース1002を具体化し得る様々な種類のシステムからデータを受信することができる。データは、場合によりヌル値を含むそれぞれのフィールド（「属性」又は「カラム」とも呼ばれる）のための値を有する記録として編成され得る。データソースからデータを最初に読み取るとき、前処理モジュール1006は、典型的には、そのデータソース内の記録に関する何らかの初期フォーマット情報から始める。一部の状況では、データソースの記録構造が最初に分からない場合があり、代わりにデータソース又はデータの分析後に決定され得る。記録に関する初期情報は、例えば個別値を表すビット数、記録内のフィールドの順序、及びビットによって表される値の種類（例えば、文字列、符号付き/符号なし整数）を含み得る。

20

【0106】

上記のモニタリング及び追跡の手法は、適切なソフトウェアを実行する計算システムを使用して実装され得る。例えば、ソフトウェアは、1つ又は複数のプログラムされた又はプログラム可能な計算システム（分散、クライアント/サーバ、グリッド等の様々なアーキテクチャのものとして実行される1つ又は複数のコンピュータプログラム内の手続きを含むことができ、かかる計算システムは、少なくとも1個のプロセッサ、少なくとも1つのデータ記憶システム（揮発性及び/又は不揮発性のメモリ及び/又は記憶素子を含む）、（少なくとも1つの入力装置又はポートを使用して入力を受け付けるための、及び少なくとも1つの出力装置又はポートを使用して出力を与えるための）少なくとも1つのユーザインタフェースをそれぞれ含む。ソフトウェアは、例えばグラフの設計、構成、及び実行に関係するサービスを提供するより大きいプログラムの1つ又は複数のモジュールを含み得る。プログラムのモジュール（例えば、グラフの要素）は、データレジスタ内に記憶されるデータモデルに適合するデータ構造又は他の編成されたデータとして実装され得る。

30

【0107】

ソフトウェアは、（例えば、汎用又は専用の計算システム又は装置によって読取可能な）CD-ROM又は他のコンピュータ可読媒体等の有形の非一時的媒体上に与えることができ、又はソフトウェアの実行場所である計算システムの有形の非一時的媒体にネットワークの通信媒体上で運ぶ（例えば、伝搬信号内に符号化する）ことができる。処理の一部又は全てを専用コンピュータ上で、又はコプロセッサ、書替え可能ゲートアレイ（FPGA）、専用の特定用途向け集積回路（ASIC）等の専用ハードウェアを使用して実行することができる。処理は、ソフトウェアによって指定される計算の様々な部分が異なる計算要素によって実行される分散方式で実装されてもよい。そのような各コンピュータプログラムは、本明細書に記載の処理を実行するために記憶装置媒体がコンピュータによって読み取られるとき、コンピュータを構成し操作するために、好ましくは汎用又は専用のプログラム可能コンピュータによってアクセス可能な記憶装置のコンピュータ可読記憶媒体（例えば、ソリッドステートメモリ若しくは媒体、又は磁気若しくは光学媒体）上に記憶

40

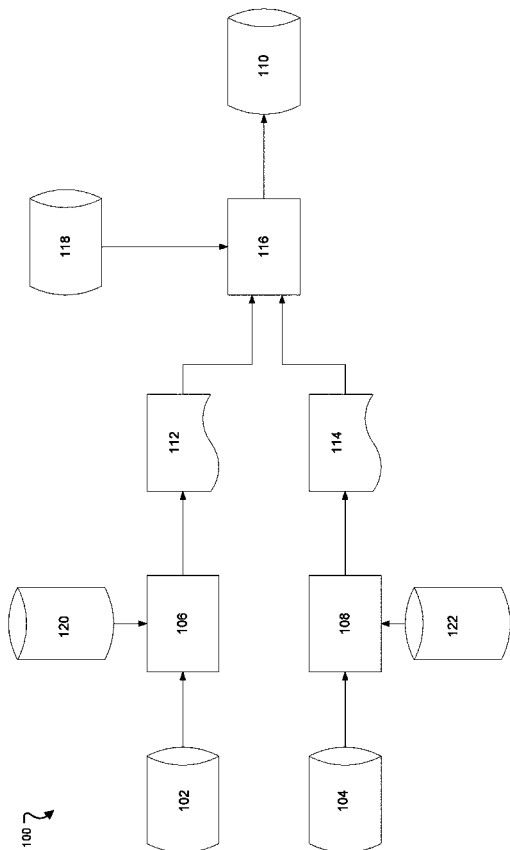
50

され又はかかるコンピュータ可読記憶媒体にダウンロードされる。本発明のシステムは、コンピュータプログラムで構成される有形の非一時的媒体として実装されると考えることもでき、そのように構成される媒体は、本明細書に記載の処理ステップの1つ又は複数を実行するためにコンピュータを特定の且つ既定の方法で動作させる。

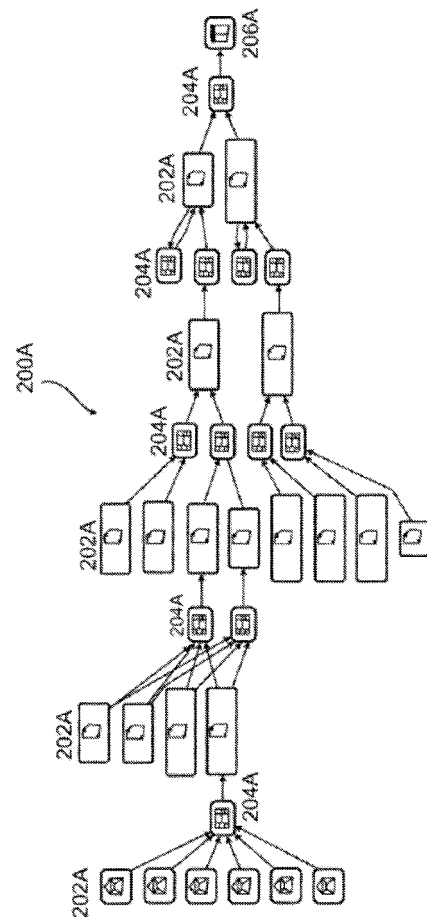
【0108】

本発明の幾つかの実施形態について説明してきた。それでもなお、上記の説明は、添付の特許請求の範囲によって定める本発明の範囲を限定するのではなく、例示を目的とすることを理解すべきである。従って、他の実施形態も添付の特許請求の範囲に含まれる。例えば、本発明の範囲から逸脱することなく様々な修正形態がなされ得る。加えて、上記のステップの一部は順序に左右されない場合があり、従って記載したのと異なる順序で実行することができる。

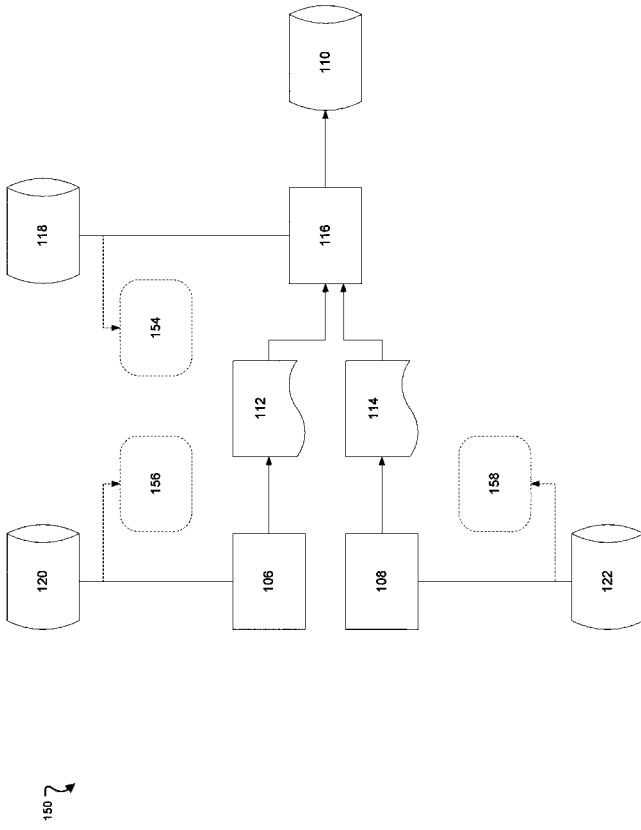
【図1】



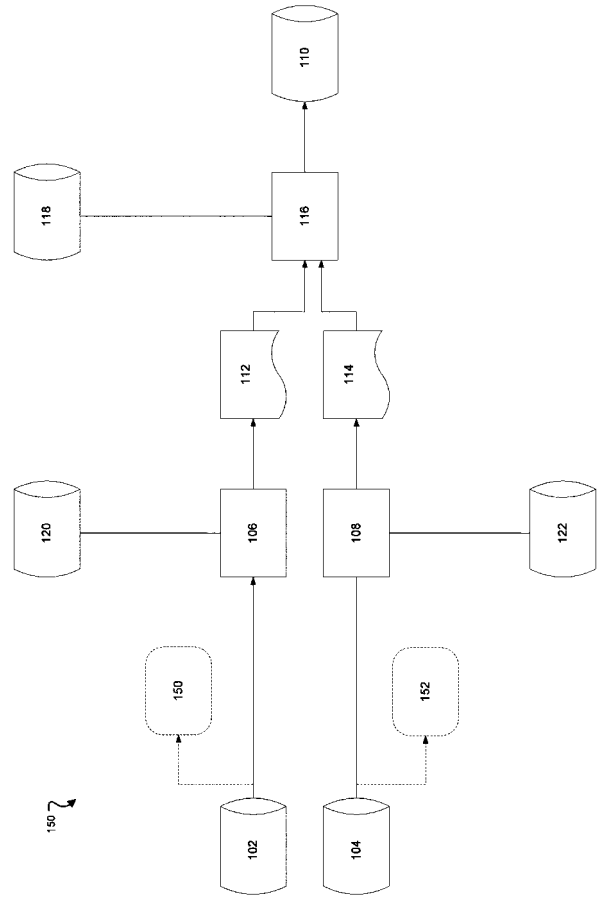
【図2】



【 図 3 A 】



【 図 3 B 】



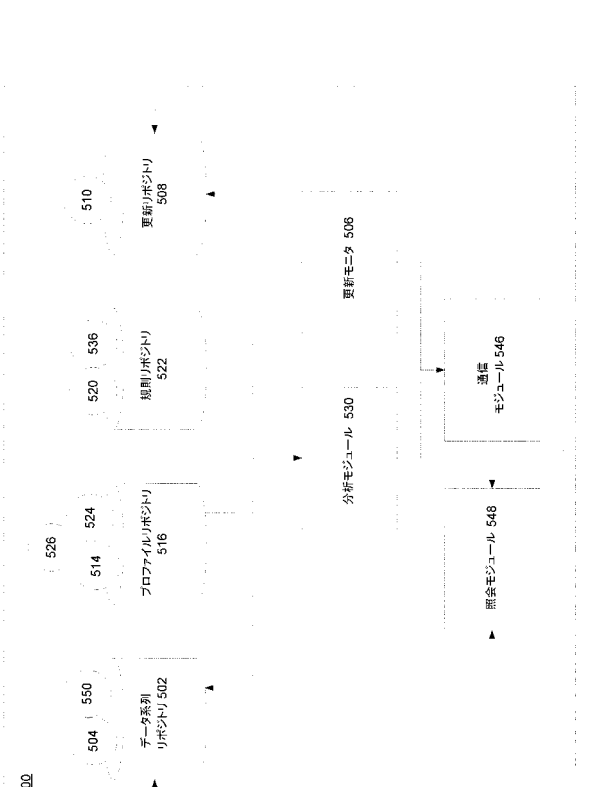
【 図 4 】

400

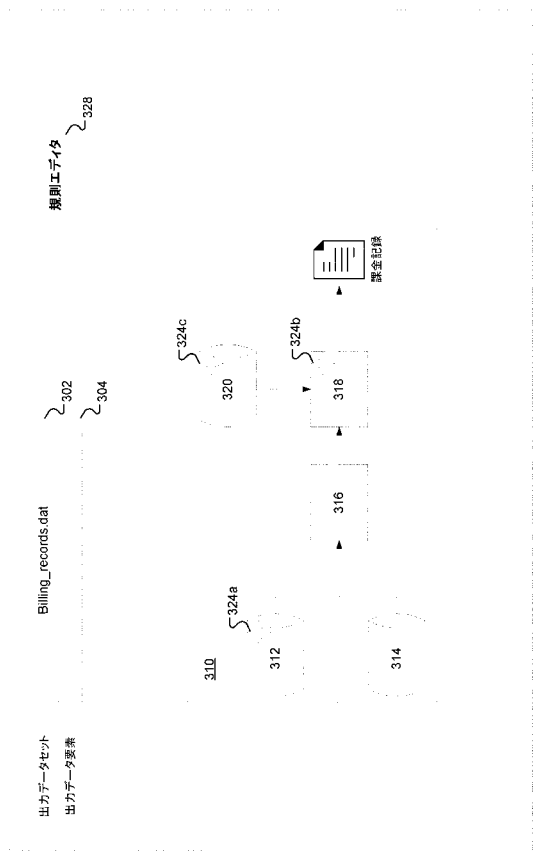
1	2	3	4	5	6	7	8	9	10	11	12	13	14	15	16	17	18	19	20	21	22	23	24	25	26	27	28	29	30	31	32	33	34	35	36	37	38	39	40	41	42	43	44	45	46	47	48	49	50	51	52	53	54	55	56	57	58	59	60	61	62	63	64	65	66	67	68	69	70	71	72	73	74	75	76	77	78	79	80	81	82	83	84	85	86	87	88	89	90	91	92	93	94	95	96	97	98	99	100
1	2	3	4	5	6	7	8	9	10	11	12	13	14	15	16	17	18	19	20	21	22	23	24	25	26	27	28	29	30	31	32	33	34	35	36	37	38	39	40	41	42	43	44	45	46	47	48	49	50	51	52	53	54	55	56	57	58	59	60	61	62	63	64	65	66	67	68	69	70	71	72	73	74	75	76	77	78	79	80	81	82	83	84	85	86	87	88	89	90	91	92	93	94	95	96	97	98	99	100

402

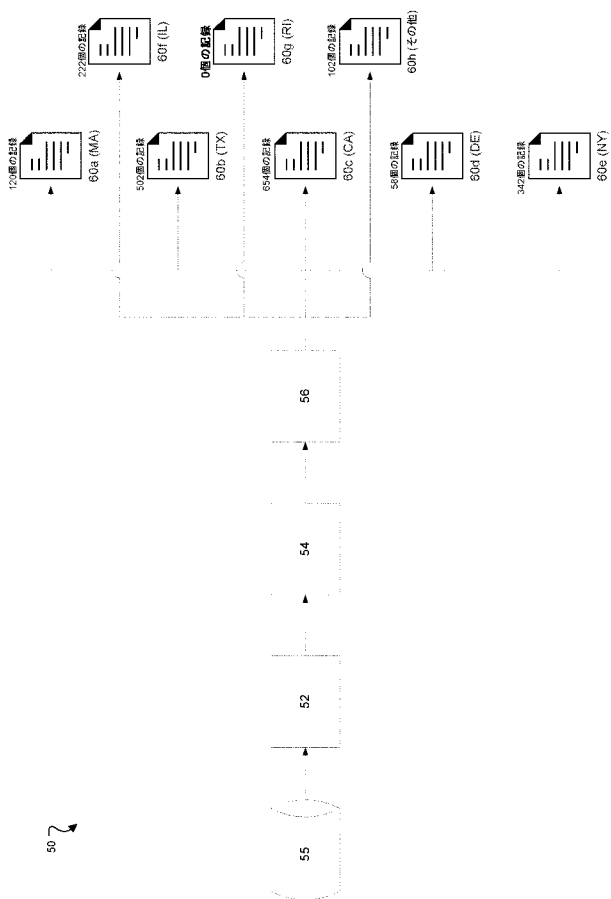
【 図 5 】



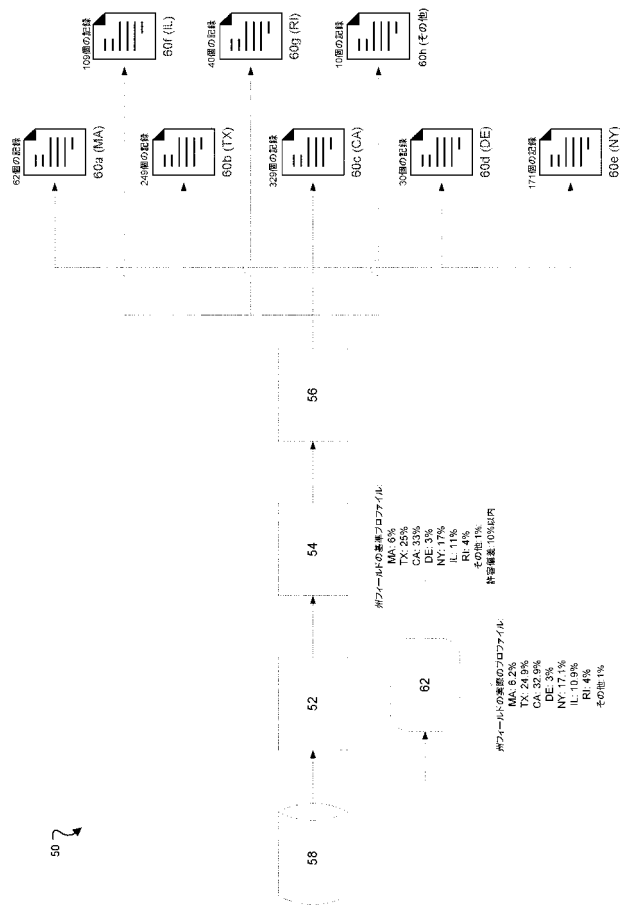
【図6】



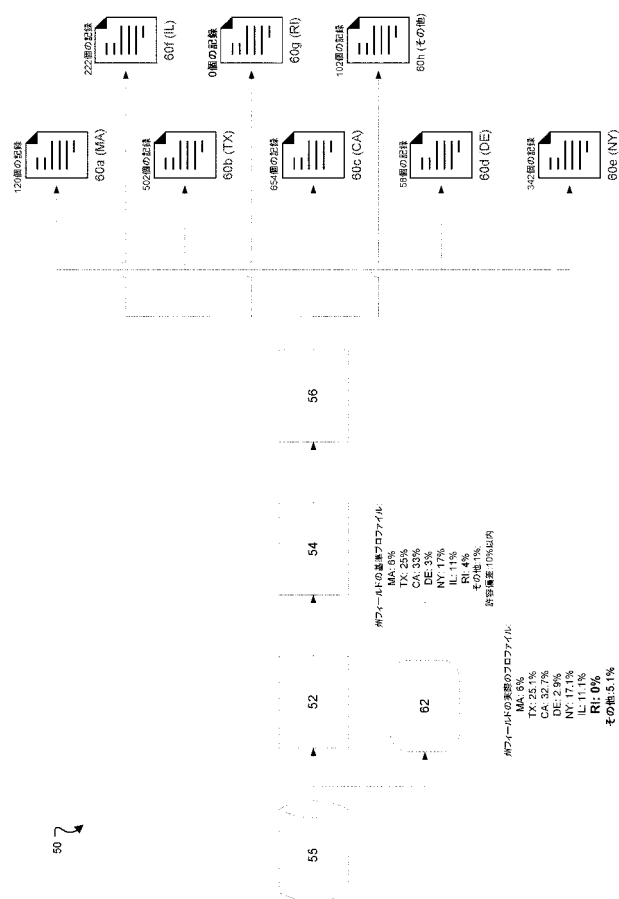
【図8A】



【図7】



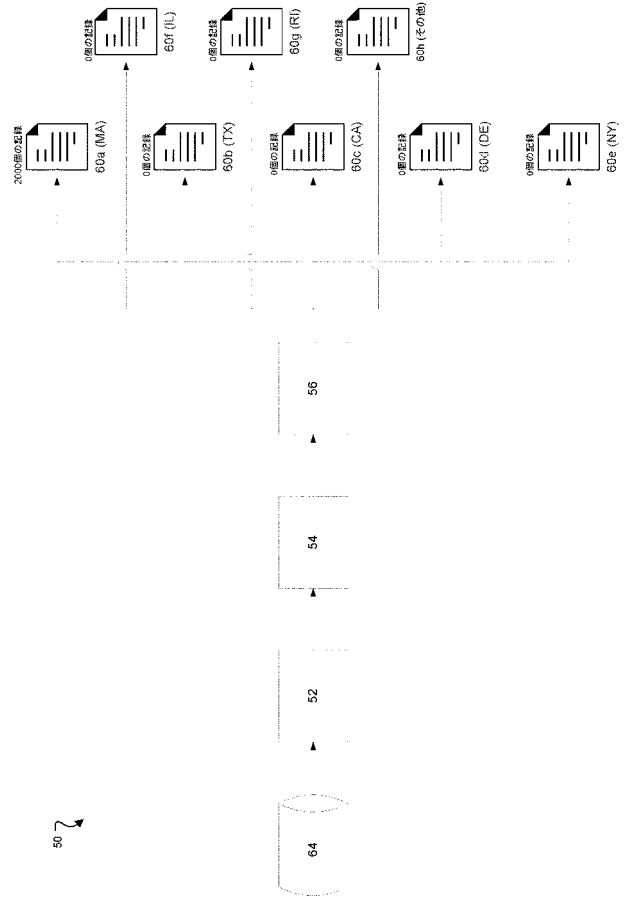
【図8B】



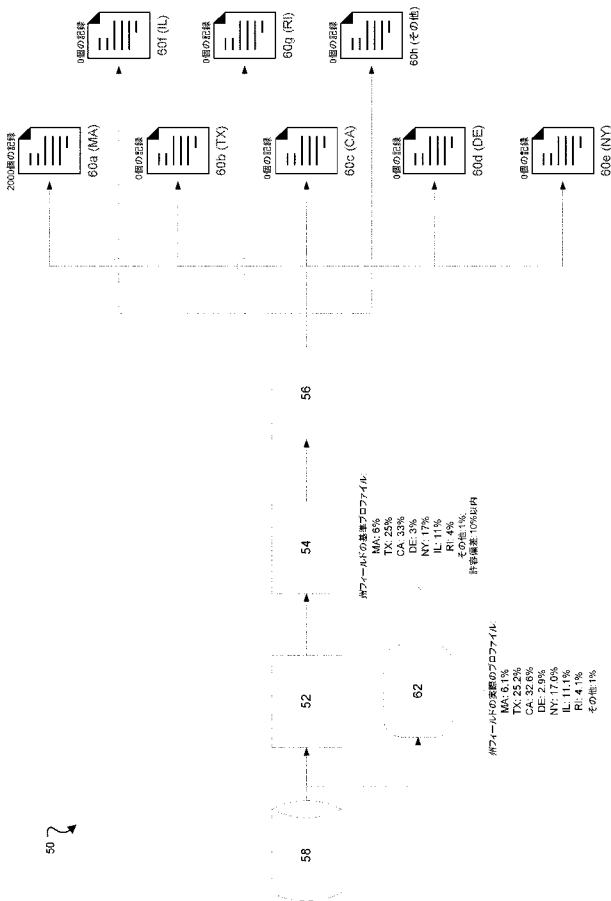
【図 8 C】

顧客_ID	州	購入_価格
12575	CA	127.51
35468	MA	547.77
26842	NY	29.80
32458	TX	245.01
54893	IL	55.76
05786	TX	18.76
65420	CA	442.10
12475	CA	54.30
54978	NY	22.99
97845	IR	65.17
65432	IL	18.19
54983	CA	25.44
32877	TX	509.60
54446	TX	63.12
23158	CA	45.18
35458	NY	56.04
87653	IR	17.10
65487	NY	152.60
32455	CA	146.12
54877	IL	33.32
78361	TX	423.01
12756	DE	18.76
45886	TX	87.45
32115	CA	108.69
68575	CA	146.08

【図 9 A】



【図 9 B】



【図 10 A】

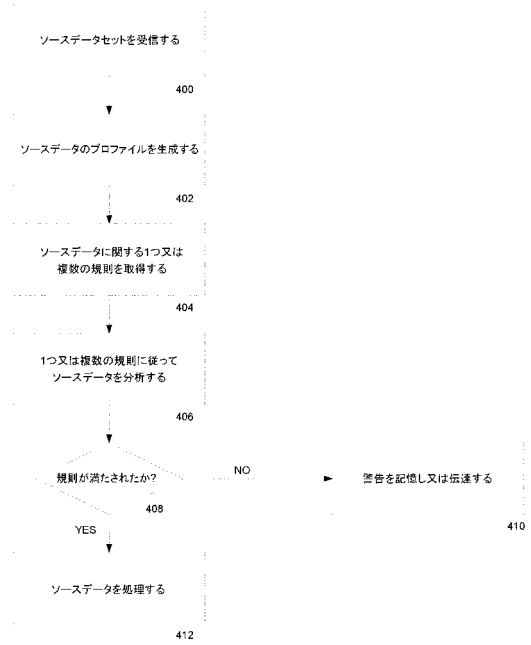


【図10B】

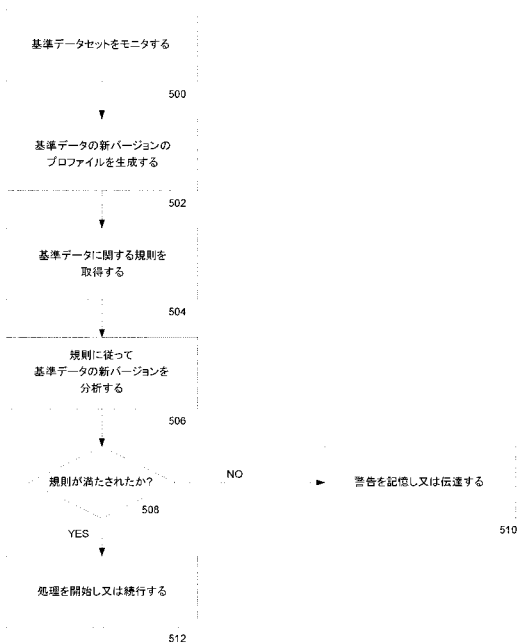
94 ↘

電話_番号	開始_時間	持続時間
6508001245	12:45:00 PM	0:30
4082604525	12:45:08 PM	1:15
212508412E	12:45:12 PM	10:12
6174438771	12:45:18 PM	8:07
3024023548	12:45:22 PM	0:02
21445893546	12:45:25 PM	12:17
9724505465	12:45:30 PM	8:32
917508459E	12:45:30 PM	0:09
7738645947	12:45:32 PM	1:12
2105782156	12:45:35 PM	4:01
4085647895	12:45:38 PM	0:14
6174546824	12:45:40 PM	2:08
214568778F	12:45:41 PM	14:16
7735684456	12:45:42 PM	9:12
6175443215	12:45:45 PM	58:17
212565788E	12:45:46 PM	7:26
3025488865	12:45:49 PM	0:16
9175453256	12:45:50 PM	1:38

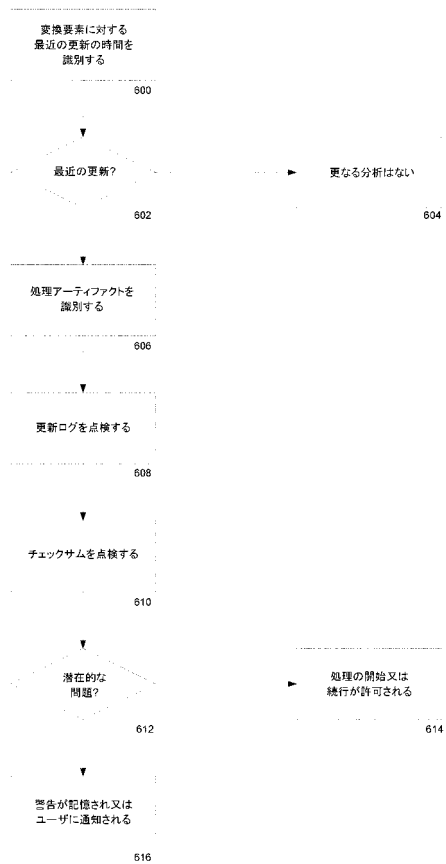
【図11】



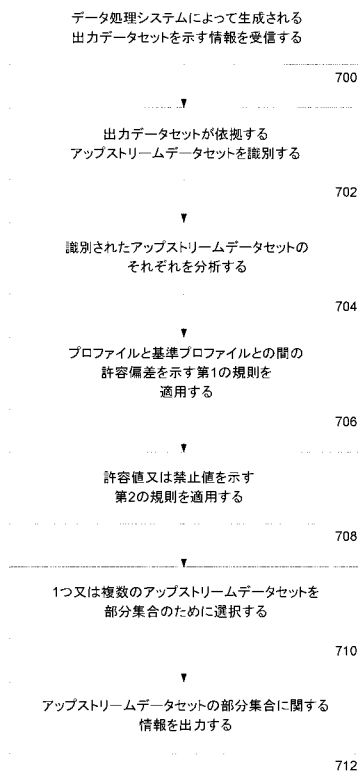
【図12】



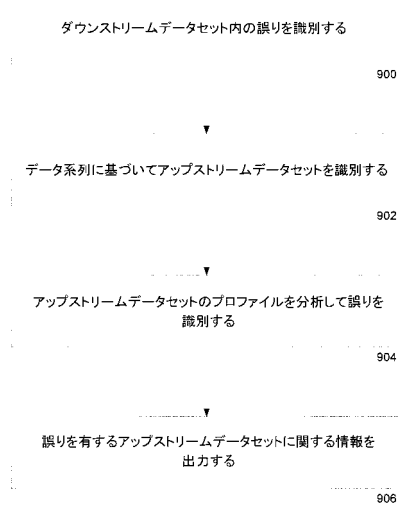
【図13】



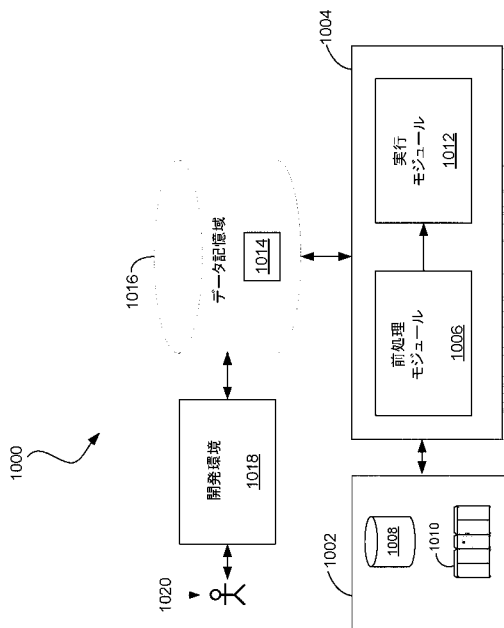
【 図 1 4 】



【 図 1 5 】



【 図 1 6 】



【手続補正書】

【提出日】令和2年6月12日(2020.6.12)

【手続補正1】

【補正対象書類名】特許請求の範囲

【補正対象項目名】全文

【補正方法】変更

【補正の内容】

【特許請求の範囲】

【請求項1】

データセットのフィールドについてのデータ品質規則を決定するためのコンピュータ実装方法であって、

前記データセットを分析することであって、前記データセットのフィールドについての基準プロファイルを決

定された前記基準プロファイルに基づいて、前記データセットの前記フィールドについてのデータ品質規則を決定することであって、前記データセットの前記フィールドについてのデータ品質規則は、

(i) 前記フィールドについての前記基準プロファイルと、前記データセットの一つ以上のデータ記録の前記フィールドについてのプロファイルとの間の許容偏差

(ii) 前記データセットのデータ記録の前記フィールドのデータ要素についての許容値、又は

(iii) 前記データセットのデータ記録の前記フィールドのデータ要素についての禁止値の一つ以上を示す、ことと、

を含むコンピュータ実装方法。

【請求項2】

請求項1に記載の方法であって、前記データセットを分析することは、前記データセットの一つ以上の履歴インスタンスを分析することを含む、方法。

【請求項3】

請求項1に記載の方法であって、前記フィールドについての前記基準プロファイルを決

定することは、前記フィールドについての履歴平均プロファイルを決

定することは、前記フィールドのデータ要素についての前記履歴平均プロファイルの変動が閾値量より少なく変動するまで前記データセットの複数の特定のインスタンスを分析することを含む、方法。

【請求項5】

請求項1に記載の方法であって、前記フィールドについての前記基準プロファイルを決

定することは、前記フィールドのデータ要素についての履歴平均値を識別することを含む、方法。

【請求項6】

請求項1に記載の方法であって、前記フィールドについての前記基準プロファイルを決

定することは、前記フィールドのデータ要素についての標準偏差値を識別することを含む、方法。

【請求項7】

請求項1に記載の方法であって、前記フィールドについての前記基準プロファイルを決

定することは、前記フィールドのデータ要素についての個別値の数を識別することを含む、方法。

【請求項8】

請求項1に記載の方法であって、前記データセットを分析することは、前記データセットの予め定められた数の特定のインスタンスを分析することを含む、方法。

【請求項9】

請求項 1 に記載の方法であって、機械学習技術を用いて、前記データセットを分析することを含む、方法。

【請求項 1 0】

請求項 1 に記載の方法であって、前記データ品質規則を前記データセットの特定のインスタンスの前記一つ以上のデータ記録に適用することを含む、方法。

【請求項 1 1】

請求項 1 0 に記載の方法であって、前記データ品質規則を前記データセットの前記特定のインスタンスに適用することは、前記データセットの前記特定のインスタンスが誤り又は起こり得る誤りを有するものと判定することを含む、方法。

【請求項 1 2】

請求項 1 1 に記載の方法であって、前記データセットの前記特定のインスタンスが誤り又は起こり得る誤りを有するものと判定することは、

前記フィールドについての前記基準プロファイルと、前記データセットの特定のインスタンスの一つ以上のデータ記録の前記フィールドについてのプロファイルとの間の偏差を判定することと、

前記基準プロファイルと前記プロファイルとの間の前記偏差が前記許容偏差を超えるものと判定することと

を含む、方法。

【請求項 1 3】

請求項 1 1 に記載の方法であって、前記データセットの前記特定のインスタンスが誤り又は起こり得る誤りを有するものと判定することは、前記許容値又は前記禁止を満たさない前記データセットの前記特定のインスタンスの一つ以上のデータ記録のフィールドのデータ要素を識別することを含む、方法。

【請求項 1 4】

データセットのフィールドについてのデータ品質規則を決定することを計算システムに行わせる命令を記憶する非一時的コンピュータ可読媒体であって、前記命令は、前記計算システムに、

前記データセットを分析することであって、前記データセットのフィールドについての基準プロファイルを決定することを含む、ことと、

決定された前記基準プロファイルに基づいて、前記データセットの前記フィールドについてのデータ品質規則を決定することであって、前記データセットの前記フィールドについてのデータ品質規則は、

(i) 前記フィールドについての前記基準プロファイルと、前記データセットの一つ以上のデータ記録の前記フィールドについてのプロファイルとの間の許容偏差

(i i) 前記データセットのデータ記録の前記フィールドのデータ要素についての許容値、又は

(i i i) 前記データセットのデータ記録の前記フィールドのデータ要素についての禁止値の一つ以上を示す、ことと、

を行わせる、非一時的コンピュータ可読媒体。

【請求項 1 5】

請求項 1 4 に記載の非一時的コンピュータ可読媒体であって、前記データセットを分析することは、前記データセットの一つ以上の履歴インスタンスを分析することを含む、非一時的コンピュータ可読媒体。

【請求項 1 6】

請求項 1 4 に記載の非一時的コンピュータ可読媒体であって、前記フィールドについての前記基準プロファイルを決定することは、前記フィールドについての履歴平均プロファイルを決定することを含む、非一時的コンピュータ可読媒体。

【請求項 1 7】

請求項 1 6 に記載の非一時的コンピュータ可読媒体であって、前記データセットを分析することは、前記フィールドのデータ要素についての前記履歴平均プロファイルの変動が

閾値量より少なく変動するまで前記データセットの複数の特定のインスタンスを分析することを含む、非一時的コンピュータ可読媒体。

【請求項 18】

請求項 14 に記載の非一時的コンピュータ可読媒体であって、前記フィールドについての前記基準プロファイルを決定することは、前記フィールドのデータ要素についての履歴平均値を識別することを含む、非一時的コンピュータ可読媒体。

【請求項 19】

請求項 14 に記載の非一時的コンピュータ可読媒体であって、前記フィールドについての前記基準プロファイルを決定することは、前記フィールドのデータ要素についての標準偏差値を識別することを含む、非一時的コンピュータ可読媒体。

【請求項 20】

請求項 14 に記載の非一時的コンピュータ可読媒体であって、前記フィールドについての前記基準プロファイルを決定することは、前記フィールドのデータ要素についての個別値の数を識別することを含む、非一時的コンピュータ可読媒体。

【請求項 21】

請求項 14 に記載の非一時的コンピュータ可読媒体であって、前記データセットを分析することは、前記データセットの予め定められた数の特定のインスタンスを分析することを含む、非一時的コンピュータ可読媒体。

【請求項 22】

請求項 14 に記載の非一時的コンピュータ可読媒体であって、前記命令は、前記計算システムに、機械学習技術を用いて、前記データセットを分析させる、非一時的コンピュータ可読媒体。

【請求項 23】

請求項 14 に記載の非一時的コンピュータ可読媒体であって、前記命令は、前記計算システムに、前記データ品質規則を前記データセットの特定のインスタンスの前記一つ以上のデータ記録に適用させる、非一時的コンピュータ可読媒体。

【請求項 24】

請求項 23 に記載の非一時的コンピュータ可読媒体であって、前記データ品質規則を前記データセットの前記特定のインスタンスに適用することは、前記データセットの前記特定のインスタンスが誤り又は起こり得る誤りを有するものと判定することを含む、非一時的コンピュータ可読媒体。

【請求項 25】

請求項 24 に記載の非一時的コンピュータ可読媒体であって、前記データセットの前記特定のインスタンスが誤り又は起こり得る誤りを有するものと判定することは、

前記フィールドについての前記基準プロファイルと、前記データセットの特定のインスタンスの一つ以上のデータ記録の前記フィールドについてのプロファイルとの間の偏差を判定することと、

前記基準プロファイルと前記プロファイルとの間の前記偏差が前記許容偏差を超えるものと判定することと

を含む、非一時的コンピュータ可読媒体。

【請求項 26】

請求項 24 に記載の非一時的コンピュータ可読媒体であって、前記データセットの前記特定のインスタンスが誤り又は起こり得る誤りを有するものと判定することは、前記許容値又は前記禁止を満たさない前記データセットの前記特定のインスタンスの一つ以上のデータ記録のフィールドのデータ要素を識別することを含む、非一時的コンピュータ可読媒体。

【請求項 27】

データセットのフィールドについてのデータ品質規則を決定するための計算システムであって、

メモリに結合されている一つ以上のプロセッサであって、前記一つ以上のプロセッサ及

び前記メモリは、

前記データセットを分析することであって、前記データセットのフィールドについての基準プロファイルを決定することを含む、ことと、

決定された前記基準プロファイルに基づいて、前記データセットの前記フィールドについてのデータ品質規則を決定することであって、前記データセットの前記フィールドについてのデータ品質規則は、

(i) 前記フィールドについての前記基準プロファイルと、前記データセットの一つ以上のデータ記録の前記フィールドについてのプロファイルとの間の許容偏差

(i i) 前記データセットのデータ記録の前記フィールドのデータ要素についての許容値、又は

(i i i) 前記データセットのデータ記録の前記フィールドのデータ要素についての禁止値の一つ以上を示す、ことと、

を行うように構成されている、計算システム。

【請求項 28】

請求項 27 に記載の計算システムであって、前記データセットを分析することは、前記データセットの一つ以上の履歴インスタンスを分析することを含む、計算システム。

【請求項 29】

請求項 27 に記載の計算システムであって、前記フィールドについての前記基準プロファイルを決定することは、前記フィールドについての履歴平均プロファイルを決定することを含む、計算システム。

【請求項 30】

請求項 29 に記載の計算システムであって、前記データセットを分析することは、前記フィールドのデータ要素についての前記履歴平均プロファイルの変動が閾値量より少なく変動するまで前記データセットの複数の特定のインスタンスを分析することを含む、計算システム。

【請求項 31】

請求項 27 に記載の計算システムであって、前記フィールドについての前記基準プロファイルを決定することは、前記フィールドのデータ要素についての履歴平均値を識別することを含む、計算システム。

【請求項 32】

請求項 27 に記載の計算システムであって、前記フィールドについての前記基準プロファイルを決定することは、前記フィールドのデータ要素についての標準偏差値を識別することを含む、計算システム。

【請求項 33】

請求項 27 に記載の計算システムであって、前記フィールドについての前記基準プロファイルを決定することは、前記フィールドのデータ要素についての個別値の数を識別することを含む、計算システム。

【請求項 34】

請求項 27 に記載の計算システムであって、前記データセットを分析することは、前記データセットの予め定められた数の特定のインスタンスを分析することを含む、計算システム。

【請求項 35】

請求項 27 に記載の計算システムであって、前記一つ以上のプロセッサ及び前記メモリは、機械学習技術を用いて、前記データセットを分析するように構成されている、計算システム。

【請求項 36】

請求項 27 に記載の計算システムであって、前記一つ以上のプロセッサ及び前記メモリは、前記データ品質規則を前記データセットの特定のインスタンスの前記一つ以上のデータ記録に適用するように構成されている、計算システム。

【請求項 37】

請求項 36 に記載の計算システムであって、前記データ品質規則を前記データセットの前記特定のインスタンスに適用することは、前記データセットの前記特定のインスタンスが誤り又は起こり得る誤りを有するものと判定することを含む、計算システム。

【請求項 38】

請求項 37 に記載の計算システムであって、前記データセットの前記特定のインスタンスが誤り又は起こり得る誤りを有するものと判定することは、

前記フィールドについての前記基準プロファイルと、前記データセットの特定のインスタンスの一つ以上のデータ記録の前記フィールドについてのプロファイルとの間の偏差を判定することと、

前記基準プロファイルと前記プロファイルとの間の前記偏差が前記許容偏差を超えるものと判定することと

を含む、計算システム。

【請求項 39】

請求項 37 に記載の計算システムであって、前記データセットの前記特定のインスタンスが誤り又は起こり得る誤りを有するものと判定することは、前記許容値又は前記禁止を満たさない前記データセットの前記特定のインスタンスの一つ以上のデータ記録のフィールドのデータ要素を識別することを含む、計算システム。

【請求項 40】

データセットのフィールドについてのデータ品質規則を決定するための計算システムであって、

前記データセットを分析する手段であって、前記データセットのフィールドについての基準プロファイルを決定することを含む、手段と、

決定された前記基準プロファイルに基づいて、前記データセットの前記フィールドについてのデータ品質規則を決定する手段であって、前記データセットの前記フィールドについてのデータ品質規則は、

(i) 前記フィールドについての前記基準プロファイルと、前記データセットの一つ以上のデータ記録の前記フィールドについてのプロファイルとの間の許容偏差

(i i) 前記データセットのデータ記録の前記フィールドのデータ要素についての許容値、又は

(i i i) 前記データセットのデータ記録の前記フィールドのデータ要素についての禁止値の一つ以上を示す、手段と、

を備える、計算システム。

フロントページの続き

1 . U N I X

2 . J A V A

(72)発明者 スピッツ, チャック

アメリカ合衆国, マサチューセッツ州 0 2 4 8 2 , ウェルズリー, アイビー ロード 3 3

(72)発明者 ゴウルド, ジョエル

アメリカ合衆国, マサチューセッツ州 0 2 4 7 4 , アーリントン, リー テラス 2 7

Fターム(参考) 5B175 FB04 KA09

【外国語明細書】
2020161147000001.pdf