

(19) 日本国特許庁 (JP)

(12) 特 許 公 報 (B2)

(11) 特許番号

特許第6843882号
(P6843882)

(45) 発行日 令和3年3月17日 (2021.3.17)

(24) 登録日 令和3年2月26日 (2021.2.26)

(51) Int. Cl.

F I

G O 6 F 16/2457 (2019.01)

G O 6 F 16/2457

請求項の数 21 (全 32 頁)

(21) 出願番号	特願2018-555888 (P2018-555888)	(73) 特許権者	515160389
(86) (22) 出願日	平成29年4月26日 (2017.4.26)		インフォマティカ エルエルシー
(65) 公表番号	特表2019-519027 (P2019-519027A)		アメリカ合衆国 カリフォルニア州 94
(43) 公表日	令和1年7月4日 (2019.7.4)		063、レッドウッド シティ、シーポー
(86) 国際出願番号	PCT/US2017/029583		ト ブルバード 2100
(87) 国際公開番号	W02017/189693		2100 Seaport Blvd, R
(87) 国際公開日	平成29年11月2日 (2017.11.2)		edwood City, CA 9406
審査請求日	令和2年4月25日 (2020.4.25)		3 U. S. A.
(31) 優先権主張番号	15/139, 186	(74) 代理人	100120662
(32) 優先日	平成28年4月26日 (2016.4.26)		弁理士 川上 桂子
(33) 優先権主張国・地域又は機関	米国 (US)	(74) 代理人	100140327
			弁理士 大塚 千秋
早期審査対象出願			
		最終頁に続く	

(54) 【発明の名称】 履歴ログからの学習と、ETLツール内のデータアセットに関するデータベースオペレーションの推奨

(57) 【特許請求の範囲】

【請求項 1】

1つ以上のコンピューティングデバイスによって、データ解析アプリケーションのインスタンスのユーザに推奨を提供する方法であって、

前記1つ以上のコンピューティングデバイスの少なくとも1つによって、コンテキストデータから、データベースオペレーション履歴データエントリおよびトレーニングコンテキストデータエントリをキャプチャすることによって、コンテキストデータのプロファイリングを行うステップであって、前記コンテキストデータは前記データ解析アプリケーション内のテーブル上で実行されるデータベースオペレーションにตอบสนองして前記データ解析アプリケーションの1つまたは複数のインスタンスから受信されたエントリを含む、ステップと、

10

前記1つ以上のコンピューティングデバイスの少なくとも1つによって、第1のセットのユーザについて、複数のテーブル上で実行される複数のデータベースオペレーションに対してプロファイリングされたデータベースオペレーション履歴データおよびプロファイリングされたコンテキストデータを維持するステップと、

前記1つ以上のコンピューティングデバイスの少なくとも1つによって、データ解析アプリケーションの第2のセットのユーザに、少なくとも1つのデータベースオペレーションまたは少なくとも1つのオペランド、の1つ以上を推薦するよう構成された複数の予測モデルを生成するステップであって、前記予測モデルのそれぞれは、プロファイルされたコンテキストデータからのコンテキストデータフィールドに対応する複数の特徴と、推薦

20

のための複数の対応するデータベースオペレーションまたは複数のオペランドのいずれかとを含む、ステップと、

アプリケーションコンテキストデータを含むアプリケーションログエントリを受信するステップであって、前記アプリケーションログエントリはデータ解析アプリケーションのインスタンス内のテーブル内の列を選択する第2のセットのユーザに応答して受信される、ステップと、

前記1つ以上のコンピューティングデバイスの少なくとも1つによって、前記アプリケーションコンテキストデータに少なくとも部分的に基づいて、前記複数の予測モデル内の1つ以上の予測モデルを選択するステップと、

前記1つ以上のコンピューティングデバイスの少なくとも1つによって、前記アプリケーションコンテキストデータを前記1つ以上の選択された予測モデルへ入力することにより、1つ以上の確率リストを生成するステップであって、前記確率リストのそれぞれは、前記複数のデータベースオペレーションまたは前記複数のオペランドに関連付けられた複数の確率値を含む、ステップと、

前記予測モデルへの入力として使用できるフォーマットでアプリケーションコンテキストデータをキャプチャするためにアプリケーションログエントリのプロファイリングを行うステップと、

前記予測モデルへの入力として前記アプリケーションコンテキストデータを使用して、1つまたは複数の推奨データベースオペレーションを決定するステップと、

前記1つ以上のコンピューティングデバイスの少なくとも1つによって、前記1つ以上の確率リストに少なくとも部分的に基づいて、1つ以上の推奨を決定するステップであって、前記1つ以上の推奨における各推奨は、データベースオペレーションまたはオペランドを含む、ステップと、

前記1つ以上のコンピューティングデバイスの少なくとも1つによって、ユーザへの提示のためにデータ解析アプリケーションのインスタンスに前記1つ以上の推奨を送信するステップと、を含む方法。

【請求項2】

前記複数の予測モデルを生成するステップが、複数の予測モデルのそれぞれについて、プロファイルされたコンテキストデータから複数のコンテキストデータフィールドを選択することにより複数の特徴を決定することと、

推薦する複数のデータベースオペレーションまたは複数のオペランドを決定することと

、
複数のデータベースオペレーションまたは複数のオペランドの各々について、複数の特徴の各々についての特徴重みを決定することとを含み、

前記特徴重みはデータベースオペレーションまたはオペランドに関する特徴の予測性の尺度に対応する、請求項1に記載の方法。

【請求項3】

前記コンテキストデータは、プロジェクトメタデータ、ワークシートメタデータ、およびユーザメタデータのうちの少なくとも1つを含む、請求項1に記載の方法。

【請求項4】

少なくとも1つの予測モデルが、多項ロジスティック分類器である、請求項1に記載の方法。

【請求項5】

前記アプリケーションコンテキストデータは、プロジェクトメタデータ、ワークシートメタデータ、およびユーザメタデータのうちの少なくとも1つを含む、請求項1に記載の方法。

【請求項6】

前記1つ以上の推奨は、ジョインオペレーションおよび結合オペレーションのうちの少なくとも1つを含む、請求項1に記載の方法。

【請求項7】

前記 1 つ以上の選択された予測モデルは、オペレーションモデルおよびオペランドモデルを含み、

前記アプリケーションコンテキストデータを前記 1 つ以上の選択された予測モデルへ入力することにより、1 つ以上の確率リストを生成するステップは、

前記アプリケーションコンテキストデータを前記オペレーションモデルへ入力することにより、前記複数のデータベースオペレーションに関連付けられた確率を含む第 1 の確率リストを生成することと、

前記アプリケーションコンテキストデータおよび前記第 1 の確率リストを前記オペランドモデルへ入力することにより、前記複数のオペランドに関連付けられた確率を含む第 2 の確率リストを生成することと、を含む、

請求項 1 に記載の方法。

【請求項 8】

データ解析アプリケーションのインスタンスのユーザに推奨を提供する装置であって、1 つ以上のプロセッサと、

前記 1 つ以上のプロセッサの少なくとも 1 つに動作可能に結合された 1 つ以上のメモリとを備え、

前記 1 つ以上のメモリは、前記 1 つ以上のプロセッサの少なくとも 1 つに実行されたときに、前記 1 つ以上のプロセッサの少なくとも 1 つに、

コンテキストデータから、データベースオペレーション履歴データエントリおよびトレーニングコンテキストデータエントリをキャプチャすることによって、コンテキストデータのプロファイリングを行うステップであって、前記コンテキストデータは前記データ解析アプリケーション内のテーブル上で実行されるデータベースオペレーションに回答して前記データ解析アプリケーションの 1 つまたは複数のインスタンスから受信されたエントリを含む、ステップと、

第 1 のセットのユーザについて、複数のテーブル上で実行される複数のデータベースオペレーションに対してプロファイリングされたデータベースオペレーション履歴データおよびプロファイリングされたコンテキストデータを維持するステップと、

データ解析アプリケーションの第 2 のセットのユーザに、少なくとも 1 つのデータベースオペレーションまたは少なくとも 1 つのオペランド、の 1 つ以上を推薦するよう構成された複数の予測モデルを生成するステップであって、前記予測モデルのそれぞれは、プロファイルされたコンテキストデータからのコンテキストデータフィールドに対応する複数の特徴と、推薦のための複数の対応するデータベースオペレーションまたは複数のオペランドのいずれかとを含む、ステップと、

アプリケーションコンテキストデータを含むアプリケーションログエントリを受信するステップであって、前記アプリケーションログエントリはデータ解析アプリケーションのインスタンス内のテーブル内の列を選択する第 2 のセットのユーザに回答して受信される、ステップと、

前記アプリケーションコンテキストデータに少なくとも部分的に基づいて、前記複数の予測モデル内の 1 つ以上の予測モデルを選択するステップと、

前記アプリケーションコンテキストデータを前記 1 つ以上の選択された予測モデルへ入力することにより、1 つ以上の確率リストを生成するステップであって、前記確率リストのそれぞれは、前記複数のデータベースオペレーションまたは前記複数のオペランドに関連付けられた複数の確率値を含む、ステップと、

前記 1 つ以上の確率リストに少なくとも部分的に基づいて、1 つ以上の推奨を決定するステップであって、前記 1 つ以上の推奨における各推奨は、データベースオペレーションまたはオペランドを含む、ステップと、

ユーザへの提示のためにデータ解析アプリケーションのインスタンスに前記 1 つ以上の推奨を送信するステップと、

を実行させる命令が格納されている、装置。

【請求項 9】

前記命令は、前記 1 つ以上のプロセッサの少なくとも 1 つに実行されたときに、前記 1 つ以上のプロセッサの少なくとも 1 つに、

前記複数の予測モデルを生成するステップにおいて、複数の予測モデルのそれぞれについて、

プロファイルされたコンテキストデータから複数のコンテキストデータフィールドを選択することにより複数の特徴を決定することと、

推薦する複数のデータベースオペレーションまたは複数のオペランドを決定することと

、
複数のデータベースオペレーションまたは複数のオペランドの各々について、複数の特徴の各々についての特徴重みを決定することと、をさらに行わせ、

前記特徴重みはデータベースオペレーションまたはオペランドに関する特徴の予測性の尺度に対応する、請求項 8 に記載の装置。

【請求項 10】

前記コンテキストデータは、プロジェクトメタデータ、ワークシートメタデータ、およびユーザメタデータのうちの少なくとも 1 つを含む、請求項 8 に記載の装置。

【請求項 11】

少なくとも 1 つの予測モデルが、多項ロジスティック分類器である、請求項 8 に記載の装置。

【請求項 12】

前記アプリケーションコンテキストデータは、プロジェクトメタデータ、ワークシートメタデータ、およびユーザメタデータのうちの少なくとも 1 つを含む、請求項 8 に記載の装置。

【請求項 13】

前記 1 つ以上の推奨は、ジョインオペレーションおよび結合オペレーションのうちの少なくとも 1 つを含む、請求項 8 に記載の装置。

【請求項 14】

前記 1 つ以上の選択された予測モデルは、オペレーションモデルおよびオペランドモデルを含み、

前記命令は、前記 1 つ以上のプロセッサの少なくとも 1 つに実行されたときに、前記 1 つ以上のプロセッサの少なくとも 1 つに、

前記アプリケーションコンテキストデータを前記 1 つ以上の選択された予測モデルへ入力することにより、1 つ以上の確率リストを生成するステップにおいて、

前記アプリケーションコンテキストデータを前記オペレーションモデルへ入力することにより、前記複数のデータベースオペレーションに関連付けられた確率を含む第 1 の確率リストを生成することと、

前記アプリケーションコンテキストデータおよび前記第 1 の確率リストを前記オペランドモデルへ入力することにより、前記複数のオペランドに関連付けられた確率を含む第 2 の確率リストを生成することと、をさらに実行させる、

請求項 8 に記載の装置。

【請求項 15】

コンピュータ可読命令を記憶する少なくとも 1 つの非一時的なコンピュータ可読記憶媒体であって、前記コンピュータ可読命令は、1 つ以上のコンピューティングデバイスに実行されたとき、前記 1 つ以上のコンピューティングデバイスの少なくとも 1 つに、

コンテキストデータから、データベースオペレーション履歴データエントリおよびトレーニングコンテキストデータエントリをキャプチャすることによって、コンテキストデータのプロファイリングを行うステップであって、前記コンテキストデータはデータ解析アプリケーション内のテーブル上で実行されるデータベースオペレーションに回答して前記データ解析アプリケーションの 1 つまたは複数のインスタンスから受信されたエントリを含む、ステップと、

第 1 のセットのユーザについて、複数のテーブル上で実行される複数のデータベースオ

10

20

30

40

50

ペレーションに対してプロファイリングされたデータベースオペレーション履歴データおよびプロファイリングされたコンテキストデータを維持するステップと、

データ解析アプリケーションの第2のセットのユーザに、少なくとも1つのデータベースオペレーションまたは少なくとも1つのオペランド、の1つ以上を推薦するよう構成された複数の予測モデルを生成するステップであって、前記予測モデルのそれぞれは、プロファイルされたコンテキストデータからのコンテキストデータフィールドに対応する複数の特徴と、推薦のための複数の対応するデータベースオペレーションまたは複数のオペランドのいずれかを含む、ステップと、

アプリケーションコンテキストデータを含むアプリケーションログエントリを受信するステップであって、前記アプリケーションログエントリはデータ解析アプリケーションのインスタンス内のテーブル内の列を選択する第2のセットのユーザに応答して受信される、ステップと、

10

前記アプリケーションコンテキストデータに少なくとも部分的に基づいて、前記複数の予測モデル内の1つ以上の予測モデルを選択するステップと、

前記アプリケーションコンテキストデータを前記1つ以上の選択された予測モデルへ入力することにより、1つ以上の確率リストを生成するステップであって、前記確率リストのそれぞれは、前記複数のデータベースオペレーションまたは前記複数のオペランドに関連付けられた複数の確率値を含む、ステップと、

前記1つ以上の確率リストに少なくとも部分的に基づいて、1つ以上の推奨を決定するステップであって、前記1つ以上の推奨における各推奨は、データベースオペレーションまたはオペランドを含む、ステップと、

20

ユーザへの提示のためにデータ解析アプリケーションのインスタンスに前記1つ以上の推奨を送信するステップと、

を実行させる、記憶媒体。

【請求項16】

前記コンピュータ可読命令は、前記1つ以上のプロセッサの少なくとも1つに実行されたときに、前記1つ以上のプロセッサの少なくとも1つに、

前記複数の予測モデルを生成するステップにおいて、複数の予測モデルのそれぞれについて、

プロファイルされたコンテキストデータから複数のコンテキストデータフィールドを選択することにより複数の特徴を決定することと、

30

推薦する複数のデータベースオペレーションまたは複数のオペランドを決定することと、

複数のデータベースオペレーションまたは複数のオペランドの各々について、複数の特徴の各々についての特徴重みを決定することと、をさらに行わせ、

前記特徴重みはデータベースオペレーションまたはオペランドに関する特徴の予測性の尺度に対応する、請求項15に記載の記憶媒体。

【請求項17】

前記コンテキストデータは、プロジェクトメタデータ、ワークシートメタデータ、およびユーザメタデータのうちの少なくとも1つを含む、請求項15に記載の記憶媒体。

40

【請求項18】

少なくとも1つの予測モデルが、多項ロジスティック分類器である、請求項15に記載の記憶媒体。

【請求項19】

前記アプリケーションコンテキストデータは、プロジェクトメタデータ、ワークシートメタデータ、およびユーザメタデータのうちの少なくとも1つを含む、請求項15に記載の記憶媒体。

【請求項20】

前記1つ以上の推奨は、ジョインオペレーションおよび結合オペレーションのうちの少なくとも1つを含む、請求項15に記載の記憶媒体。

50

【請求項 2 1】

前記 1 つ以上の選択された予測モデルは、オペレーションモデルおよびオペランドモデルを含み、

前記コンピュータ可読命令は、前記 1 つ以上のプロセッサの少なくとも 1 つに実行されたときに、前記 1 つ以上のプロセッサの少なくとも 1 つに、

前記アプリケーションコンテキストデータを前記 1 つ以上の選択された予測モデルへ入力することにより、1 つ以上の確率リストを生成するステップにおいて、

前記アプリケーションコンテキストデータを前記オペレーションモデルへ入力することにより、前記複数のデータベースオペレーションに関連付けられた確率を含む第 1 の確率リストを生成することと、

前記アプリケーションコンテキストデータおよび前記第 1 の確率リストを前記オペランドモデルへ入力することにより、前記複数のオペランドに関連付けられた確率を含む第 2 の確率リストを生成することと、をさらに実行させる、

請求項 1 5 に記載の記憶媒体。

【発明の詳細な説明】**【技術分野】****【0001】**

本出願は 2016 年 4 月 26 日に出願された米国非仮出願第 15 / 139 , 186 号の優先権を主張し、その開示は、その全体が参照により本明細書に組み込まれる。

【0002】

本開示は一般に、データベース管理システムおよびデータウェアハウスにおけるデータプロセスの抽出、変換、およびロードに関し、より詳細には、データ閲覧および編集環境において表示されるデータのためのデータベースオペレーションを決定し、推奨するためのコンピュータ実行方法に関する。

【背景技術】**【0003】**

データウェアハウスの分野では、複数の外部データソースからのデータが通常、内部データベース管理システムに取り込まれるときに、抽出 (extract)、変換 (transform)、およびロード (load) (ETL) プロセスを介して遷移する。ETL プロセスの一部として、データは、(i) 1 つまたは複数のデータソースから抽出され、(ii) 内部データソースのビジネス要件および技術要件に従ってプログラム変換され、(iii) 内部データベース管理システムのターゲットデータストアにロードされる。一旦システムに入ると、データは、様々なデータベースオペレーションを使用してシステムユーザによって操作され得る。多くの場合、ユーザは膨大な量のデータを扱っており、一部のユーザは、データベース管理アプリケーションがデータを処理するためにサポートするデータベースオペレーションに慣れていないか、またはデータベース管理システム内でデータを処理する最も効率的な方法を知らない。この問題に対処するのに十分な知識および経験を獲得することは、特に、一時ユーザまたは多くのタイプのデータを扱うユーザにとって困難であり、時間がかかる可能性がある。

【発明の概要】**【0004】**

データ解析サーバは、機械学習予測モデルを使用して、プログラマ的に決定された推奨データベース動作を、データ解析アプリケーションの習熟度が低いユーザ (ガイド付きユーザ) に提供するように構成される。予測モデルは、データベース内の類似データに関する上級ユーザ (トレーニングユーザ) によるデータベースオペレーション入力から学習される。予測モデルは、熟練度の低いユーザがどのデータベースオペレーションがデータに適しているかを選択するプロセスを改善することによって、データベースを操作する際の効率を改善することを可能にする。

【0005】

データ解析サーバは、以前のデータベースユーザによるデータベースオペレーションの

10

20

30

40

50

履歴データを使用して、ETLツールのユーザにデータベースオペレーションを推奨するための予測モデルを構築する。データプロファイリングモジュールは、選択されたユーザグループに提示され、ユーザによって操作されるデータベーステーブルおよびテーブルセット（プロジェクト）のコンテキストデータを維持するように構成される。コンテキストデータは、テーブルおよびプロジェクトのメタデータを含む。データベースオペレーション履歴モジュールは、テーブルおよびプロジェクト上のデータベースオペレーションの履歴データを維持するように構成される。本明細書で使用されるデータベースオペレーションは、ETLによってサポートされ、変換または変更されたデータセットを生成するために特定のデータに対して実行されるプログラム操作である。特定のデータベースオペレーションには、ジョイン（コンバイン）、結合（マージ）、フィルタ、フォーミュラ、ルックアップ、列分割、列追加（データ拡張）、パターン認識および不整合修正、データクレンジング、データ整合、データ標準化などが含まれる。データベース演算は、数学的演算、方程式などのデータに対する演算をさらに含むことができる。

10

【0006】

データベースオペレーション推奨モジュールはデータベースオペレーションをユーザに推奨するための予測モデルを構築し、トレーニングし、使用するように構成される。データベースオペレーション推奨モジュールは維持されたデータベースオペレーション履歴データおよびコンテキストデータを使用してモデルをトレーニングし、それによって、どのコンテキストデータが特定のデータベースオペレーションの適用を予測するかを決定する。ガイド付きユーザによるデータベースの使用中に、リアルタイムでガイド付きユーザに対する推薦を生成するために、データベースオペレーション推薦モジュールはガイド付きユーザによってアクセスされている特定のテーブルまたはプロジェクトに対するコンテキストデータを受信し、予測モデルを使用してそのテーブルまたはプロジェクトに対して実行する1つまたは複数の推薦データベースオペレーションを決定する。

20

【0007】

データ解析アプリケーションのグラフィカルユーザインタフェースは、データセクション、情報セクション、および様々なユーザインタフェース制御を含む。データセクションは、分析用のテーブルを表示するためのものである。情報セクションは、テーブルのスキーマ定義に基づいて、テーブルのプロファイル情報を表示するものである。コンポジットデータ制御はテーブル間の少なくとも1つのマッチング列に基づいてテーブルをコンポジットテーブルに統合するデータベースオペレーション（同等には、データベースコマンド）を受け取るためのものである。複合データ制御は、様々な統一データベースオペレーションのための複数の異なる制御であってもよい。UIの推薦制御は、データベースオペレーション推薦モジュールによって決定された推薦データベースオペレーションを表示するためのものである。

30

【0008】

本明細書に記載される特徴および利点はすべてを包含するものではなく、特に、多くの追加の特徴および利点が、図面、明細書、および特許請求の範囲を考慮して、当業者には明らかであろう。さらに、本明細書で使用される言語は主に、読みやすさおよび説明の目的のために選択されており、本発明の主題を描写または限定するために選択されていない場合があることに留意されたい。

40

【図面の簡単な説明】

【0009】

【図1】図1は、一実施形態による、データ解析アプリケーションにおいて、データベースオペレーションの履歴ログから予測モデルを生成し、データに対するデータベースオペレーションを推奨するコンピューティング環境の高レベルブロック図である。

【図2】図2は、一実施形態によるデータベースオペレーション推薦モジュールのより詳細な図を示す。

【図3】図3は、予測モデルをトレーニングする際に使用するための特徴およびクラスを示す例示的なデータテーブルである。

50

【図 4】図 4 は、一実施形態による、データ解析アプリケーションにおいてデータを閲覧および操作するためのユーザインタフェースの一例を示す。

【図 5 A】図 5 A は一実施形態による、データ解析アプリケーションのガイド付きユーザに対してデータベースオペレーションを決定し推奨するための予測モデルを構築し、トレーニングするための方法を示すフローチャートである。

【図 5 B】図 5 B は、一実施形態による、データ解析アプリケーションのガイド付きユーザにデータベースオペレーションを推薦するためにトレーニングされた予測モデルを使用する方法を示すフローチャートである。

【図 6】図 6 は、一実施形態による、選択された列に応答して提供される推奨を備えた、図 3 の例示的なユーザインタフェースを示す。

10

【図 7】図 7 は、データ解析アプリケーションにおいて、データ解析サーバから受信した推奨データベースオペレーションおよびオペランドを提示するための方法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】

【0010】

システムのアーキテクチャ

図 1 は、一実施形態による、データ解析アプリケーションにおいて、データベースオペレーションの履歴ログから予測モデルを生成し、データに対するデータベースオペレーションを推奨するコンピューティング環境 100 の高レベルブロック図である。

【0011】

20

示されるように、コンピューティング環境 100 は、データリポジトリ 102、データ解析サーバ 104、およびデータ解析アプリケーション 125 を含む。

【0012】

複数のデータリポジトリ 102（本明細書では個別にデータリポジトリ 102 と呼ぶ）は、データを管理するための 1 つまたは複数のシステムを含む。各データリポジトリ 102 は、データリポジトリ 102 内に格納されたデータにアクセスして更新するためのチャネルを提供する。データリポジトリ 102 内のデータは、ユーザ、ユーザのグループ、エンティティ、および/またはワークフローに関連付けられ得る。例えば、データリポジトリ 102 は、特定のエンティティに関連付けられたすべての個人に関連付けられたデータを記憶する顧客関係管理（CRM）システムまたは人事（HR）管理システムとすることができる。データリポジトリ 102 は、ETL プロセスのためのデータソースまたはエクスポートターゲットとすることができる。データソースの例は、データベース、アプリケーション、およびローカルファイルを含む。同様に、これらのソースは、データをエクスポートするためのターゲットとして機能することができる。共通のエクスポートターゲットは、TABLEAU、SALESFORCE WAVE、およびEXCELである。

30

【0013】

データ解析アプリケーション 125 は、ユーザがデータ解析サーバ 104 によってデータリポジトリ 102 から抽出されたデータを操作し、単一のテーブル又は多数のテーブルに対して実行されるべきデータベースオペレーションを選択及び指定することを可能にするソフトウェアアプリケーションであり、この機能を実行するための 1 つの手段である。一実施形態では、データ解析アプリケーション 125 がテーブルのセットであるプロジェクトの形でユーザにデータを提供する。データ解析アプリケーション 125 の様々なモジュールは、汎用コンピュータシステムのネイティブコンポーネントまたは標準コンポーネントではなく、コンピュータシステムの汎用機能を超えて拡張する、本明細書で説明する特定の機能を提供する。さらに、モジュールの機能および動作はコンピュータシステムによる実装を必要とするほど十分に複雑であり、したがって、いかなる実際的な実施形態でも、人間の心の中の精神的なステップによって実行することはできない。これらの構成要素の各々は、以下により詳細に記載される。データ解析アプリケーション 125 はデバイス非依存であり、したがって、デスクトップアプリケーション、モバイルアプリケーション、またはウェブベースのアプリケーションとすることができる。その様々な機能を実行

40

50

するために、データ解析アプリケーション 125 は、ユーザインタフェース (UI) モジュール 122 およびデータベースオペレーション UI モジュール 124 を含む。

【0014】

いくつかの実施形態では、データ解析アプリケーション 125 は、様々なオンサイトおよび外部のソースおよびターゲット、ならびに本明細書で説明されるプロセスに關与する強化サービスと共に、より大きなクラウドアーキテクチャの一部である。

【0015】

UI モジュール 122 は UI において表示するためのデータを受信し、受信したデータに対応するユーザインタフェースを生成し、受信したデータをテーブルにポピュレートし、予測モデルに基づいてデータリファインメントの推奨を表示し、テーブルの 1 つまたは複数の列に關連付けられた列サマリを生成し、これらの機能を実行するための 1 つの手段である。生成されたユーザインタフェースは、データ解析アプリケーション 125 のユーザがテーブルエントリを操作すること、およびデータベースオペレーションをデータに適用することを含めて、テーブルを見ること、およびテーブルと対話することを可能にする。

10

【0016】

データベースオペレーション UI モジュール 124 は UI モジュール 122 によって生成されたテーブル内のデータに適用するための 1 つ以上のデータベースオペレーション制御を提供し、この機能を実行するための 1 つの手段である。具体的には、データベースオペレーション UI モジュール 124 がデータ解析アプリケーション 125 のユーザがテーブルに關連付けられたデータベースオペレーションを選択し、指定し、および / またはデータベースオペレーションの適用を引き起こすことを可能にする制御を提供する。

20

【0017】

一実施形態によれば、UI モジュール 122 およびデータベースオペレーション UI モジュール 124 によって提供されるユーザインタフェースは、グラフィカルに表現されたデータセクション、情報セクション、および様々なグラフィカルに表現されたデータベースオペレーション制御を含む。UI のデータ部は、解析用のテーブルを表示するためのものである。UI の情報セクションは、テーブルに關するプロファイル情報を表示するためのものである。プロファイル情報は、コンテキストデータなどのテーブルの特徴を記述する。UI の複合データ制御は、テーブル間の少なくとも 1 つの一致する列に基づいて 2 つのテーブルを複合テーブルに統合するコマンドを受信するユーザインタフェース要素である。UI の推薦制御は、予測モデルを用いてデータベースオペレーション推薦モジュール 114 により決定された推薦データベースオペレーションを表示するユーザインタフェース要素である。UI は、図 4 および図 6 に關して以下により詳細に説明される。

30

【0018】

データベースオペレーション UI モジュール 124 は、実行された各データベースオペレーションに対して、表示されたテーブルに対して実行された各データベースオペレーションをデータ解析サーバ 104 内のデータベースオペレーション履歴モジュール 112 に送信する。各データベースオペレーションは、オペレーション識別子によって表され、オペレーション識別子は例えば、名前、ID 番号、およびデータベースオペレーションに含まれていたオペランドを示すオペレーション記述によって、オペレーションを一意に識別する。データベースオペレーション履歴モジュール 112 は、データに適用されたデータベースオペレーションをデータベースオペレーション履歴記憶部 120 に記憶する。経時的にデータに適用されるデータベースオペレーションはデータベースオペレーション履歴ストア 120 に取り込まれ、データベースオペレーション履歴内の任意のステップはアンドゥ、リドゥ、または異なるデータに適用することができる。データベースオペレーションは後述するように、ログの形式で格納することができる。

40

【0019】

データ解析サーバ 104 はデータをデータリポジトリ 102 から抽出し、データを処理し、処理されたデータをデータ解析アプリケーション 125 に提供して、データをユーザ

50

に表示し、ユーザによって操作できるようにする。これらの機能を実行するために、データ解析サーバ104は、データ抽出モジュール108と、データプロファイリングモジュール110と、データベースオペレーション履歴モジュール112とを含む。さらに、これらの機能に関連するデータを記憶するために、データプロファイリングサーバ104は、リポジトリデータストア116、プロファイリングデータストア118、およびデータベースオペレーション履歴ストア120を含む。分析サーバ104の様々なモジュールは、汎用コンピュータシステムのネイティブコンポーネントまたは標準コンポーネントではなく、コンピュータシステムの汎用機能を超えて拡張する、本明細書で説明する特定の機能を提供する。さらに、モジュールの機能および動作はコンピュータシステムによる実装を必要とするほど十分に複雑であり、したがって、いかなる実際的な実施形態でも、人間の心の中の精神的なステップによって実行することはできない。これらの構成要素の各々は、以下により詳細に記載される。

10

【0020】

データ抽出モジュール108は抽出されるべきデータリポジトリ102内のデータを識別し、そのデータをデータリポジトリ102から取り出し、そのデータをリポジトリデータストア116に格納するように構成され、そのための1つの手段である。動作中、データ抽出モジュール108は、データを抽出する1つまたは複数のデータリポジトリ102を識別する。データ抽出モジュール108はまた、抽出されるべき識別されたデータリポジトリ102に記憶された特定のデータを識別する。データリポジトリ102および/またはそこに格納された特定のデータの識別は、データプロファイリング動作を行うユーザから受け取った命令に基づいて行うことができる。あるいは、そのような識別がデータを抽出する外部データソースを指定する1つまたは複数のビジネスロジック定義に基づいて行うことができる。

20

【0021】

データ抽出モジュール108は、データリポジトリ102によって提供されるデータアクセスチャネルを介してデータリポジトリ102から識別されたデータを抽出する。一実施形態では、データ・アクセス・チャネルは、データ抽出モジュール108がデータ・リポジトリ102と安全に通信して、データ・リポジトリ102との間でデータを取り出し、送信することを可能にする安全なデータ転送プロトコルである。データがデータリポジトリ102から抽出されると、データ抽出モジュール108は、データをリポジトリデータストア116に格納する。

30

【0022】

データプロファイリングモジュール110はデータリポジトリ102から抽出され、リポジトリデータストア116に格納されたデータを処理して、データのすべての列、行、および領域を完全にプロファイリングし、そうするための1つの手段である。列、行、およびデータフィールドのプロファイリングは、データタイプ、データドメイン、およびエントリ長、固有値パーセント、および空白値パーセントなどのデータ値に関する他の情報を識別することを含む。

【0023】

データベース運用履歴部112は、セル、テーブル、プロジェクトに適用されるデータベース運用の履歴を受け取り、格納する手段の一つである。動作中、データベース動作がセル、テーブル、またはプロジェクトに適用されるとき、データベース動作履歴モジュール112は、適用された特定のデータベース動作と、どのデータに適用されたかを、データベース動作履歴ストア120に記憶する。したがって、経時的にデータに適用されるデータベースオペレーションは、データベースオペレーション履歴ストア120に取り込まれる。

40

【0024】

本明細書で使用されるデータベースオペレーションは、ETLシステムのプログラムコードによってサポートされ、変換または変更されたデータセットを生成するために特定のデータに対して実行されるプログラム操作である。データベースオペレーションは、テ

50

ブルまたはプロジェクトに対して実行することができる。特定のデータベースオペレーションには、ジョイン（コンバイン）、結合（マージ）、フィルタ、フォーミュラ、ルックアップ、列分割、列追加（データ拡張）、パターン認識および不整合修正、データクレンジング、データ整合、データ標準化が含まれる。データベース演算は、数学的演算、数式などのデータに対する演算をさらに含むことができる。データベースオペレーションには以下のものがある。

【 0 0 2 5 】

【 表 1 】

表 1：データベースオペレーション

データベース オペレーション	説明
ジョイン (join)	2つのソース間の列の1つまたは複数のペアと一致する条件に基づいてソースを結合
結合 (union)	異なるソースのデータを同じフィールド名でマージ
expr	オペランドとして与えられる数式を評価
フィルタ (filter)	テーブルまたは列の行を指定された条件に基づいてフィルタリング
ルックアップ (lookup)	別のデータベーステーブルまたはビューの値を返す
列追加 (column add)	新しい列を追加（行ごとに）。新しい列は、空、デフォルト値、または式に基づいてポピュレートできる。
列分割 (column split)	カラムを複数のカラムに分割。得られる各列は、元の列の一部を有する。
グループ分け (group by)	指定した列の固有値に基づいて行をグループ化
集約 (aggregate)	行のグループ上で式を計算
ソート (sort)	指定した条件に基づいて行または列を並べ替え

【 0 0 2 6 】

データベースオペレーション履歴モジュール 1 1 2 はさらに、抽出されたデータに関するコンテキストデータを受信し、作成し、管理するように構成される。コンテキストデータは、テーブルまたはプロジェクトに対して実行されているデータベースオペレーションに関連して収集または生成されるテーブルおよび / またはプロジェクトに関する情報である。コンテキストデータは、プロジェクトメタデータ、テーブルメタデータ、列メタデータ、およびユーザメタデータを含む。コンテキストデータは、データベースオペレーション履歴ストア 1 2 0 に格納されてもよい。

【 0 0 2 7 】

プロジェクトメタデータフィールドは、以下を含む。

10

20

30

40

【表 2】

表 2 : プロジェクトメタデータフィールド

メタデータフィールド	説明	変数タイプ
project_id	プロジェクト ID	英数字/数字の固有 ID
project_name	プロジェクト名	テキスト文字列
num_worksheets	プロジェクト内のテーブル数	整数
num_join_worksheets	プロジェクト内でジョインされたテーブル数	整数
num_union_worksheets	プロジェクト内で結合されたテーブル数	整数
num_agg_worksheets	プロジェクト内で実行された集合オペレーションの数	整数

10

【 0 0 2 8 】

テーブルメタデータフィールドは以下を含む。

【表 3】

表 3 : テーブルメタデータフィールド

20

メタデータフィールド	説明	変数タイプ
ws_id	テーブル ID	英数字/数字の固有 ID
ws_name	テーブル名	テキスト文字列
ws_type	ワークシートタイプ	値に基づくテーブルタイプ
ws_rows	テーブル内の行数	整数
ws_curr_size	テーブルサイズ (単位?)	整数
ws_unique_cols	テーブル内の固有列数	整数
ws_text_cols	テキストを含む列の数	整数
ws_date_cols	日付を含む列の数	整数
ws_numeric_cols	数値を含む列の数	整数 (whole number)
ws_blank_cols	空白列の数	整数 (whole number)
ws_hidden_cols	隠された列の数	整数 (whole number)
ws_derived_cols	派生値を含む列の数	整数 (whole number)
recipe	テーブル上で実行されたオペレーションのシーケンス	テキスト文字列のリスト

30

40

【 0 0 2 9 】

列メタデータフィールドフィールドは以下を含む。

【表 4】

表 4：列メタデータフィールド

メタデータフィールド	説明	変数タイプ
column_id	列 ID	英数字の固有 ID
column_name	列名	テキスト文字列
column_datatype	列内のデータタイプ	
column_nulls	列内のヌル値パーセンテージ	小数
column_unique	列内の固有値パーセンテージ	小数
column_trimmable	列内のトリミング可能な値のパーセンテージ	小数
column_outlier	列内の外れ値のパーセンテージ	小数
column_pattern	カラム値のパターン	小数
column_domain	列のドメイン	テキスト文字列
column_selection	選択された列の値	テキスト文字列
column_maxvalue	列内の最大値	整数
column_minvalue	列内の最小値	整数

10

20

【0030】

一実施形態では、コンテキストデータは、プロジェクトメタデータ、テーブルメタデータ、列メタデータ、ユーザメタデータ、および操作を含むログファイルに含まれる。ログファイル内のログエントリは、テーブルまたはプロジェクト上で実行されているデータベースオペレーションに回答して生成され、JavaScript Object Notation (JSON) で表すことができる。ログエントリは、コンテキストおよび操作履歴データを以下の形式で表現する。

30

```
{<user metadata><project metadata><worksheets metadata><column metadata><operation specifics>}
```

【0031】

ログ・エントリは、データベースオペレーション履歴記憶装置 120 に記憶することができる。ログエントリの例を以下に示す：

【0032】

ログエントリ例 1

【数 1】

```

{
  "type": "com.informatica.dataprep.suggestion.logger.ContextCollectorImpl",
  "data": {
    "userLogger": {
      "type": "com.informatica.dataprep.suggestion.logger.UserContextImpl",
      "data": {
        "user_id": 197,
        "licence_plan": "NA"
      }
    },
    "projectLogger": {
      "type": "com.informatica.dataprep.suggestion.logger.ProjectContextImpl",
      "data": {
        "project_id": 2312,
        "project_name": "test-log",
        "num_worksheets": 1,
        "num_join_worksheets": 0,
        "num_union_worksheets": 0,
        "num_agg_worksheets": 0
      }
    },
    "sheetLogger": [
      {
        "type": "com.informatica.dataprep.suggestion.logger.SheetContextImpl",
        "data": {
          "ws_id": 2313,
          "ws_name": "dp_user_session.csv",
          "ws_type": "NORMAL",
          "ws_rows": 31275,
          "ws_curr_size": 6,
          "ws_unique_cols": 3,
          "ws_text_cols": 3,
          "ws_date_cols": 0,
          "ws_numeric_cols": 3,
          "ws_blank_cols": 0,
          "ws_hidden_cols": 0,
          "ws_derived_cols": 0,
          "recipe": "deleteHeaderRows;"
        }
      }
    ],
    "columnLogger": [
      {
        "type":
"com.informatica.dataprep.suggestion.logger.ColumnContextImpl",
        "data": {
          "column_id": 2327,
          "column_name": "D",
          "column_datatype": "Integer",
          "column_nulls": 0.0,
          "column_unique": 99.81,
          "column_trimmable": 0.0,
          "column_outlier": 36.17585931254996,
          "column_pattern": "<NUMBER>",
          "column_domain": "None",
          "column_selection": "None",
          "column_maxvalue": "1427703590101",
          "column_minvalue": "1403021779000"
        }
      }
    ]
  }
}

```

【数 2】

```

"operationLogger": {
  "type":
"com.informatica.dataprep.suggestion.logger.OperationContextImpl",
  "data": {
    "operation": "expr:",
    "operation_description": "expr:(((C3/60)/60)/24000)+DATE(1970,1,1)"
  },
  "timestamp": 1427863549753
}
}

```

10

【0033】

ログエントリ例1では、ユーザメタデータが「userLogger」セクションに含まれる。このセクションの「type」サブセクションは、データがユーザメタデータ（「...UserContextImpl」）であることを示す。このセクションの「data」サブセクションは、データベースオペレーションが実行されたときにデータ解析アプリケーション125のユーザを一意に識別するユーザ識別子値（「user__id」：197）を含む。

【0034】

プロジェクトメタデータは「projectLogger」セクションに含まれる。このセクションの「type」サブセクションは、データがプロジェクトコンテキストデータ（「...ProjectContextImpl」）であることを示す。「data」サブセクションは、データベースオペレーションが実行されたプロジェクトの特徴を含む。この特徴は、プロジェクト識別子（「project__id」：2312）、プロジェクト名（「project__name」：「test-log」）、プロジェクト内のテーブルの数（「num__worksheets」：1）、ジョインされたワークシートの数（「num__join__worksheetsL」：0）、接合されたワークシートの数（「num__union__worksheets」：0）、および集合ワークシートの数（「num__agg__worksheets」：0）を含む。

20

【0035】

テーブルメタデータは、「sheetLogger」セクションに含まれる。このセクションの「type」サブセクションは、データがテーブルメタデータ（「SheetContextImpl」）であることを示す。「data」サブセクションは、データベースオペレーションが実行されたテーブルの特性を含む。この特性は、テーブル識別子（「ws__id」：2313）、テーブル名（「ws__name」：「dp__user__session.csv」）、テーブルタイプ（「ws__type」：「NORMAL」）、テーブル内の行数（「ws__rows」：31275）、テーブルサイズ（「ws__curr__size」：6）、テーブル内の固有の列数（「ws__unique__cols」：3）、テーブル内のテキストの列数（「ws__text__cols」：3）、日付形式の列数（「ws__date__cols」：0）、数の列数（「ws__numeric__cols」：3）、空白の列数（「ws__blank__cols」：0）、隠れ列数（「ws__hidden__cols」：0）、派生列数（「ws__derived__cols」：0）、およびテーブル上で実行される操作のリスト（「recipe」：「deleteHeaderRows;」）を含む。

30

40

【0036】

列メタデータは「columnLogger」セクションに含まれる。このセクションの「type」サブセクションは、データが列メタデータ（「...ColumnContextImpl」）であることを示す。「data」サブセクションは、データベースオペレーションが実行された列の特性を含む。この特性は、列識別子（「column__id」：2327）、列名（「column__name」：「D」）、列データ型（「c

50

column_datatype」:「Integer」)、列内のヌル値のパーセンテージ(「column_nulls」:0.0)、列内の一意の値のパーセンテージ(「column_unique」:99.81)、トリミングのパーセンテージ、列内の可能な値(「column_trimmable」:0.0)、列内の外れ値のパーセンテージ(「column_outlier」:36.175859312554996)、列の値のパターン(「column_pattern」:「<NUMBER>」)、列のドメイン(「column_domain」:「None」)、列の選択された領域(「column_selection」:「None」)、列の最大値(「column_maxvalue」:「1427703590101」)、および列の最小値(「column_minvalue」:「14030217779000」)を含む。

10

【0037】

データベースオペレーション履歴データは、「operationLogger」セクションに含まれる。このセクションの「type」サブセクションは、データがオペレーション履歴データ(「OperationContextImpl」)であることを示す。このセクションの「data」サブセクションは、どのデータベースオペレーションが実行されたかを識別するオペレーション識別子(「operation」:「expr:」)と、データベースオペレーションに含まれたオペランドを示すオペレーション記述(「operation_description」:「expr(((C3/60)/60)/24000)+DATE(1970,1,1)」)とを含む。この例では、オペレーションが、タイムスタンプミリ秒を日数に変換するために使用され、それらを日付1/1/1970に追加して、タイムスタンプの日付を取得する。

20

【0038】

ログエントリ例2

【数 3】

```

{
  "type": "com.informatica.dataprep.suggestion.logger.ContextCollectorImpl",
  "data": {
    "userLogger": {
      "type": "com.informatica.dataprep.suggestion.logger.UserContextImpl",
      "data": {
        "user_id": 352,
        "licence_plan": "NA"
      }
    },
    "projectLogger": {
      "type": "com.informatica.dataprep.suggestion.logger.ProjectContextImpl",
      "data": {
        "project_id": 688,
        "project_name": "DP Tables",
        "num_worksheets": 2,
        "num_join_worksheets": 0,
        "num_union_worksheets": 0,
        "num_agg_worksheets": 0
      }
    },
    "sheetLogger": [
      {
        "type": "com.informatica.dataprep.suggestion.logger.SheetContextImpl",
        "data": {
          "ws_id": 762,
          "ws_name": "dp_user.csv",
          "ws_type": "NORMAL",
          "ws_rows": 5420,
          "ws_curr_size": 7,
          "ws_unique_cols": 2,
          "ws_text_cols": 6,
          "ws_date_cols": 0,
          "ws_numeric_cols": 1,
          "ws_blank_cols": 0,
          "ws_hidden_cols": 0,
          "ws_derived_cols": 0,
          "recipe": "deleteHeaderRows;upper;"
        }
      },
      {
        "type": "com.informatica.dataprep.suggestion.logger.SheetContextImpl",
        "data": {
          "ws_id": 689,
          "ws_name": "dp_user_session.csv",
          "ws_type": "NORMAL",

```

【数 4】

```

        "ws_rows": 31275,
        "ws_curr_size": 8,
        "ws_unique_cols": 3,
        "ws_text_cols": 3,
        "ws_date_cols": 1,
        "ws_numeric_cols": 3,
        "ws_blank_cols": 1,
        "ws_hidden_cols": 0,
        "ws_derived_cols": 0,
        "recipe": "deleteHeaderRows;expr
(((C3/60)/60)/24000)+DATE(1970,1,1);expr;"
    }
},
"columnLogger": [
    {
        "type":
"com.informatica.dataprep.suggestion.logger.ColumnContextImpl",
        "data": {
            "column_id": 802,
            "column_name": "id",
            "column_datatype": "Integer",
            "column_nulls": 0.0,
            "column_unique": 100.0,
            "column_trimmable": 0.0,
            "column_outlier": 35.283262594574644,
            "column_pattern": "<NUMBER>",
            "column_domain": "None",
            "column_selection": "None",
            "column_maxvalue": "5947716",
            "column_minvalue": "22"
        }
    },
    {
        "type":
"com.informatica.dataprep.suggestion.logger.ColumnContextImpl",
        "data": {
            "column_id": 794,
            "column_name": "USER_ID",
            "column_datatype": "Integer",
            "column_nulls": 0.0,
            "column_unique": 11.55,
            "column_trimmable": 0.0,
            "column_outlier": 11.210590266675194,
            "column_pattern": "<NUMBER>",
            "column_domain": "None",
            "column_selection": "None",
            "column_maxvalue": "5947257",
            "column_minvalue": "22"
        }
    }
],
"operationLogger": {
    "type":
"com.informatica.dataprep.suggestion.logger.OperationContextImpl",
    "data": {
        "operation": "Join",
        "operation_description": "Join:762 FULL OUTER 689"
    }
},
"timestamp": 1427888644562
}

```

【0039】

ログエントリ 1 の例は、1 つのテーブルに対して実行されたデータベースオペレーションに対応する。ログエントリ 2 の例は、2 つのテーブルに対して実行されたデータベースオペレーションに対応する。2 つのテーブルに対して実行されるデータベースオペレーションは、2 つのテーブルからの列を結合するジョイン (join) および結合 (union)

10

20

30

40

50

n) 操作を含む。ログエントリ 2 の例では、Sheet Logger セクションで指定されているように、テーブル ID 762 および 689 を有するテーブルに対して完全外部ジョイン操作が実行された。ログエントリ 2 の例は、2 組のテーブルデータおよび 2 組の列データを有し、各組は、ジョイン操作が実行された 2 つのテーブルのうちの 1 つに対応する。

【0040】

データ解析サーバ 104 のユーザモジュール 115 は、ユーザがデータ解析サーバ 104 とのアカウントを管理することを可能にする。ユーザモジュール 115 はさらに、データ解析アプリケーション 125 に関連するユーザ活動に対応するユーザ情報を受信し、記憶する。ユーザ情報はユーザの好み、ユーザに関連するコンピューティングデバイスに関する情報、様々なグループ（例えば、企業（enterprise）、組織（organization）など）とのユーザの関連、およびトレーニングユーザおよび/またはガイド付きユーザとしてのユーザのステータスを含み得る。トレーニングユーザはデータ解析アプリケーション 125 のユーザであり、そのデータベースオペレーションは、ガイド付きユーザにデータベースオペレーションを推薦するための予測モデルをトレーニングするために使用される。ガイド付きユーザは、トレーニングされた予測モデルからデータベースオペレーションの推奨を受信するデータ解析アプリケーション 125 のユーザである。ガイド付きユーザの 1 つ以上のセットはトレーニングユーザに関連するデータを使用してガイド付きユーザの推薦が生成されるように、トレーニングユーザの 1 つ以上のセットに関連付けられてもよい。

【0041】

ガイド付きユーザおよび/またはトレーニングユーザとしてのユーザのステータス、ならびにガイド付きユーザのセットとトレーニングユーザとの間の関連付けは、システム管理者、他のユーザによって指定されてもよいし、自動的に指定されてもよい。例えば、グループ（例えば、組織または企業）は、データ解析アプリケーション 125 の上級ユーザをトレーニングユーザとして、経験の少ないユーザをガイド付きユーザとして指定することができる。トレーニングユーザのセットはまた、地理的領域またはデータ解析アプリケーション 125 による熟練度の尺度などのユーザ特性に基づいて、ユーザモジュール 115 によって自動的に決定されてもよい。ガイド付きユーザは、トレーニングユーザに関連付けられたトレーニングデータがガイド付きユーザに対する推薦を生成するために使用されるように、トレーニングユーザに関連付けられてもよい。結果として、トレーニングユーザの知識および経験は、データ解析サーバ 104 によって活用されて、ガイド付きユーザに有用な推薦を提供し得る。グループからのトレーニングユーザのセットを同じグループからのガイド付きユーザのセットに関連付けることにより、システムはユーザに、そのグループに特に関連する推薦を提供することができ、グループ内のユーザがグループ全体の一貫性を維持し、独自の情報（例えば、方程式、関数、およびデータ）を保護しながら、生産性を高めることができるようになる。

【0042】

一実施形態では、トレーニングユーザおよびガイド付きユーザの複数のセットが存在する。特定のユーザは同時にトレーニングユーザおよびガイド付きユーザとすることができ、複数の組のトレーニングユーザおよび/またはガイド付きユーザに属することができる。ユーザはあるタイプのプロジェクト（例えば、会計）に関してはトレーニングユーザであってもよいが、別のタイプのプロジェクト（例えば、マーケティング）に関してはガイド付きユーザであってもよい。トレーニングユーザおよび/またはガイド付きユーザとしてのユーザのステータス、ならびにトレーニングユーザとガイド付きユーザとの間の任意の関連付けは、ユーザデータストア 117 に格納され得る。ユーザモジュール 115 は、特定のプロジェクトについて、ユーザのステータスを、ガイド付きユーザまたはトレーニングユーザのいずれかとして決定することができる。ユーザがガイド付きユーザである場合、ユーザモジュール 115 はさらに、推薦を生成するためにトレーニングユーザのどのセットが使用されるべきかを決定することができる。

【 0 0 4 3 】

データベースオペレーション推薦モジュール 1 1 4 は、コンテキストデータおよびデータベースオペレーション履歴データに基づいて、ユーザに対して推薦されるデータベースオペレーションを決定する。データベースオペレーション推薦モジュール 1 1 4 は、予測モデルに基づいてデータベースオペレーションを推薦する。データベースオペレーションは、予測モデルによっても決定されるオペランドを含む。オペランドは、関数入力などのデータベースオペレーションのための入力またはパラメータである。様々な実施形態では、予測モデルがデータベース動作履歴データおよびコンテキストデータを使用することによってトレーニングされ得る機械学習アルゴリズムである。ロジスティック回帰、ニューラルネットワーク、決定木モデル、およびサポートベクトルマシンモデルを含む、様々な予測モデルが当技術分野で周知である。モデルは入力の特定のセット（例えば、コンテキストデータ）が与えられると、特定のデータベースオペレーションが適切である確率を予測し、可能性のあるデータベースオペレーションのうちの 1 つまたは複数、および任意選択で、推奨された操作に対応するオペランドを推奨する。予測モデルは、データベース動作履歴データおよびコンテキストデータを使用してトレーニングされる機械学習アルゴリズムとすることができる。一実施形態では、多項ロジスティック分類器または他の適切な汎用機械学習技法などの識別モデルが使用される。方程式、パラメータ、および他のモデル特性は、データベースオペレーション推奨ストア 1 2 1 に格納され得る。データベースオペレーション推奨を生成するための 3 つのモデル例について、図 2 を参照して以下に説明する。

【 0 0 4 4 】

図 2 は、一実施形態によるデータベースオペレーション推薦モジュール 1 1 4 のより詳細な図を示す。モデル構築モジュール 2 0 5 は予測モデルを構築し、モデルトレーニングモジュール 2 1 0 はトレーニングユーザからのトレーニングデータを使用して予測モデルをトレーニングし、推薦生成モジュール 2 2 0 は、トレーニングされた予測モデルを使用して、ガイド付きユーザに対する推薦のためのデータベース動作を決定する。一実施形態では、モデルが多項ロジスティック分類子を使用する。ログエントリからプロファイルされたメタデータフィールドによって表されるような特定のコンテキストデータが与えられると、多項式ロジスティック分類子を使用するモデルは、それぞれの確率を有するデータベース演算のリストを生成する。モデルは、トレーニングデータを用いてトレーニングされる。一実施形態では、トレーニングデータがトレーニングユーザのセットに関する格納されたデータベースオペレーション履歴データおよびコンテキストデータを含む。この実施形態では、モデルトレーニングモジュール 2 1 0 が例えばユーザデータストア 1 1 7 からモデルのトレーニングユーザを決定し、データベース動作履歴ストア 1 2 0 からトレーニングデータを取り出す。

【 0 0 4 5 】

モデル構築モジュール 2 0 5 は予測モデルを構築し、この機能を実行するための 1 つの手段である。多項ロジスティック分類器は、所与の情報に基づいて事象が発生する確率の推定値を提供する。多項ロジスティック分類器は、以下の形式をとる：

【 0 0 4 6 】

【 数 5 】

$$P(c|d) = \frac{\exp(\sum_i (\lambda_{i,c} F_{i,c}(d,c)))}{\sum_c \exp(\sum_i (\lambda_{i,c} F_{i,c}(d,c)))}$$

【 0 0 4 7 】

ここで、 $P(c|d)$ は、特徴 F によって特徴付けられる条件 d が与えられた場合に生じるクラス c によって特徴付けられる事象の確率の推定値である。クラス c は、演算またはオペランドのいずれかである特定の予測モデルの出力に対応し、特徴 F は、関連するコ

ンテキストデータに対応する。 $F_i(d, c)$ は特徴*i*の観測の尺度であり、*F*値が高いほど、特徴の存在の相対的な尺度が高いことを示す。 $w_{i,c}$ は、クラス*c*に対応する特徴*i*の特徴重みである。特定の特徴に対する高い $w_{i,c}$ は、*F*値がクラス*c*に対する強力な指標であることを示す。特徴は、異なるクラス*c*に対して異なる*F*値または $w_{i,c}$ 値を有することができる。 $P(c|d)$ によって表される確率は、クラス*c*について、クラスの全ての特徴にわたる観測の尺度と特徴の重みとの積の合計の指数(exponential)を決定し、その値を全てのクラスにわたる同じ値の合計で割ることによって、計算される。

【0048】

一実施形態では、モデル構築モジュール205が3つのモデルを構築する：演算モデル(OPモデル)、オペランドモデル(OPDモデル)、列演算モデル(OPCモデル)である。3つのモデルの各々は、トレーニングユーザーからのトレーニングデータを使用してモデルトレーニングモジュール210によってトレーニングされる。3つのモデルの各々は、推薦生成モジュール220によって使用されて、推薦されたデータベースオペレーションおよび/またはオペランドと、コンテキストデータに基づく関連する相対確率とのリストを生成する。

【0049】

OPモデルは、単一テーブルデータベースオペレーションの推奨データベースオペレーションのリストおよび関連する確率を生成する。OPモデルの機能は、列メタデータフィールドである。

【0050】

OPDモデルは、単一テーブルデータベースオペレーションのための推奨データベースオペレーションのためのオペランドのリストおよび関連する確率を生成する。OPDモデルの特徴は、列メタデータフィールドおよびデータベースオペレーションである。一実施形態では、OPDモデルがOPモデルと併せて使用され、OPモデルによって決定されたデータベースオペレーションのためのオペランドを決定する。OPDモデルは、OPDモデルによって決定された推奨オペランドが決定された演算に対応するように、OPモデルによって決定されたデータベース演算を入力として取り込む。

【0051】

OPCモデルは、2テーブルデータベースオペレーションのための推奨データベースオペレーションのリストおよび関連する確率を生成する。OPCモデルの特徴は、2つのテーブルの各々および2つの列の各々についてのメタデータである。

【表5】

表5：モデルクラスと特徴

モデル	クラス	特徴
OP	データベースオペレーション	列メタデータ
OPD	オペランド	列メタデータ、データベースオペレーション
OPC	ジョイン、結合	テーブルメタデータ、列メタデータ

【0052】

各モデルについて、モデルトレーニングモジュール210は、多項ロジスティック分類器に含める特徴としてどのコンテキストデータフィールドが選択されるかを決定する。モデルトレーニングモジュール210はさらに、選択された各特徴に対する特徴重みを決定する。すべてのメタデータフィールドが演算および/またはオペランドを予測するわけではないので、すべてのメタデータフィールドがモデルの特徴として使用されるわけではな

い。一実施形態では、モデルトレーニングモジュール210が複数のデータベース動作履歴エントリにわたって、取られる特定のデータベース動作または使用されるオペランドを予測するモデル特徴として使用するコンテキストデータフィールドを選択する。モデルトレーニングモジュール210は、各コンテキストデータフィールドについて予測性の尺度を計算し、予測性の尺度は例えば、情報利得であってもよい。各クラスについて、モデルトレーニングモジュール210は、格納されたコンテキストデータに基づいて、可能な特徴のリスト内の各特徴についての情報利得を計算する。モデルトレーニングモジュール210は、閾値情報ゲイン値を超える特徴を選択し、モデルに含める。所与のクラスについて、特徴に関する情報利得は、以下の式によって計算することができる。

$$IG(C|F) = Entropy(C) - Entropy(C|F)$$

10

ここで、 $IG(C|F)$ は情報利得であり、 $Entropy(C)$ はクラスCのエントロピーであり、 $Entropy(C|F)$ は特徴の存在を仮定したクラスCの条件付きエントロピーである。

【0053】

一実施形態では、モデルトレーニングモジュール210が情報利得を計算する前に、コンテキストデータを前処理する。一実施形態では、モデルトレーニングモジュール210がコンテキストデータを再サンプリングして、各クラスにわたるデータエントリの分布をより均一にし、その結果、より少ない頻度のデータベースオペレーションがモデルにおいて過少に表されないようにする。再サンプリング技術は、アンダーサンプリング法、オーバーサンプリング法、またはハイブリッド法を含むことができる。一実施形態では、リサンプリングがSMOTE(Synthetic Minority Oversampling Technique)を用いて実行される。様々な実施形態では、すべてのデータを数値表現に変換すること、データの正規化、および数値の2進数への量子化など、他の前処理ステップがコンテキストデータに対して実行される。

20

【0054】

図3は、OPモデルなどの予測モデルをトレーニングする際に使用するための特徴およびクラスを示すデータエントリの例示的なテーブルである。図3の例では、列301~308に示される特徴が、データベースオペレーションが実行された列に対応する列メタデータエントリの選択されたセットである。

【0055】

30

列301は、表4で識別される「column__id」メタデータフィールドからの値を含む。

【0056】

列302は、表4で識別される「column__type」メタデータフィールドからの値を含む。

【0057】

列303は、表4で識別される「column__nulls」メタデータフィールドからの値を含む。

【0058】

列304は、表4で識別される「column__unique」メタデータフィールドからの値を含む。

40

【0059】

列305は、表4で識別される「Column__pattern」メタデータフィールドからの値を含む。

【0060】

列306は、表4で識別される「column__domain」メタデータフィールドからの値を含む。

【0061】

列307は、表4で識別される「column__maxvalue」メタデータフィールドからの値を含む。

50

【 0 0 6 2 】

列 3 0 8 は、表 4 で識別される「column_min_value」メタデータフィールドからの値を含む。

【 0 0 6 3 】

ここでのモデルのクラスは列 3 1 0 に示される表 1 に識別されるように、列上で実行されたデータベースオペレーションの名前である。表 5 に示されるように、これらの特定の例示的な特徴およびクラスは、OP モデルをトレーニングするために使用される。図 3 の例は 1 4 個のデータエントリを示すが、実際には上述した予測モデルが数百、数千、数百万またはそれ以上のデータエントリを用いてトレーニングすることができる。様々な実施形態において、データエントリを構成するコンテキストデータおよびデータベースオペレーション履歴データの断片は、図 2 に関して上述したように、モデルトレーニングモジュール 2 1 0 によってログエントリから選択される。データエントリは、データベースオペレーション推奨ストア 1 2 1 に格納されてもよい。

10

【 0 0 6 4 】

推奨生成モジュール 2 2 0 は、トレーニングされた予測モデルを使用して、ガイド付きユーザに対する推奨のためのデータベースオペレーションおよび / またはオペランドのリストを、それぞれの相対確率と共に決定する。推奨生成モジュール 2 2 0 は、例えばログファイルの形式でコンテキストデータを受信する。推奨生成モジュール 2 2 0 は関連コンテキストデータを、予測モデルに入力され得るフォーマットでキャプチャするように、ログファイルをプロファイルする。推奨生成モジュール 2 2 0 は推奨を生成するために、適切な予測モデルにコンテキストデータを入力する。様々な実施形態では、使用される予測モデルが単一テーブル推奨の場合には OP モデルおよび OPD モデルであり、マルチテーブル推奨の場合には OPC モデルである。推奨生成モジュール 2 2 0 は、様々なイベントの発生時に、定期的な間隔で、または任意の他の適切な時間に、推奨を生成することができる。一実施形態では、推奨生成モジュール 2 2 0 がデータ解析アプリケーション 1 2 5 のユーザインタフェースにおける列の選択を検出し、それに応答してその列に対する推奨を生成するプログラムコードを実行する。このような推奨を生成するためのプロセスは、図 5 に関して以下に説明される。

20

【 0 0 6 5 】

推奨生成モジュール 2 2 0 は、生成された 1 つまたは複数のリストから 1 つまたは複数の推奨データベースオペレーションおよび / またはオペランドを選択する。一実施形態では、推奨生成モジュール 2 2 0 が予測モデルによって計算されるように、最も高い相対確率を有する推奨を選択する。例えば、選択された列に対する単一シート推奨に対して、推奨生成モジュール 2 2 0 は、OP モデルによって決定される 3 つの最も確からしいデータベースオペレーションと、OPD モデルによって決定される各オペレーションに対する 1 つの最も確からしいオペランドとを選択することができる。

30

【 0 0 6 6 】

推奨生成モジュール 2 2 0 はユーザに表示するために、データ解析アプリケーション 1 2 5 に推奨を提供する。一実施形態では、推奨が動作のテキスト記述として提供される。各データベースオペレーションのテキスト記述は、データベースオペレーション推奨ストア 1 2 1 に格納することができる。推奨生成モジュール 2 2 0 はユーザに表示するためにデータ解析アプリケーション 1 2 5 に提供するために、推奨されたデータベースオペレーションのためのテキスト記述を取り出してもよい。

40

【 0 0 6 7 】

図 4 は、一実施形態による、データ解析アプリケーションにおいてデータを閲覧および操作するためのユーザインタフェース 4 0 0 の一例を示す。例示的なユーザインタフェースは、データセクション 4 1 0、情報セクション 4 1 5、およびコントロール 4 1 7 を含む。

【 0 0 6 8 】

データセクション 4 1 0 は、閲覧および操作のためのテーブルを表示する。データセク

50

ション 4 1 0 は 1 つ以上のデータソース（例えば、1 0 2）から抽出されたデータでポピュレート（populated）される。この例では、2 つのテーブルタブ 4 0 5 が示され、「MDM 顧客データ（MDM Customer Data）」と題するテーブルがデータセクション 4 1 0 に表示される。ユーザは、テーブルタブ 4 0 5 を使用してプロジェクト内の他のテーブルにナビゲートすることができる。図 4 の例では、列「first_name」4 0 7 が選択される。

【 0 0 6 9 】

情報セクション 4 1 5 は、テーブルおよび選択されたデータに関するプロファイル情報を表示する。情報セクション 4 1 5 において、オーバービューカード 4 2 0 は選択された列（first_name）の情報オーバービュー（例えば、タイプ、固有値のパーセンテージ、ブランク値のパーセンテージ、列内の名前の最小長、列内の名前の最大長、およびドメインの数）を提供する。ドメインカード 4 2 5 は、テーブル 4 0 5 内のすべてのドメインに関する情報、およびどのくらいの行が各ドメインに対応するかについての情報を含む。値頻度カード 4 3 0 は、選択されたファーストネーム列 4 0 7 における種々の名前の値の頻度、並びに名前の各時間がどのように発生するかをリストする。

10

【 0 0 7 0 】

提案カード 4 3 5 は、データベースオペレーション推奨モジュール 1 1 4 によって決定された推奨データベースオペレーションを実行するための提案をユーザに提供する。図示の例では、提案されたデータベース動作がファーストネーム（first_name）として検証される。システムは、インタフェースのユーザにこれらのインテリジェントな提案を提供するのを助けるために、上述のデータプロファイリングを使用する。提案カード 4 3 5 については、図 5 および図 6 に関して以下でより詳細に説明する。

20

【 0 0 7 1 】

コントロール 3 1 7 は、ユーザが表示されたデータおよびテーブルを操作することを可能にし、データおよびテーブルに対してデータベースオペレーションを実行することを含む。データおよびテーブルは、データエントリとの対話（セル内容の編集、セルの右クリック、方程式の挿入など）または提案カード 3 3 5 などの情報セクション内の要素との対話などの他の方法で操作することもできる。

【 0 0 7 2 】

図 5 A は一実施形態による、データ解析アプリケーションのガイド付きユーザに対してデータベースオペレーションを決定し推奨するための予測モデルを構築し、トレーニングするための方法を示すフローチャートである。データ解析サーバ 1 0 4 は、データ解析アプリケーション 1 2 5 のユーザをトレーニングするためのコンテキストデータおよびデータベースオペレーション履歴データを維持する（5 0 0）。データ解析サーバ 1 0 4 はある期間にわたって、データ解析アプリケーション 1 2 5 のインスタンスからのコンテキストデータおよびデータベースオペレーション履歴データを、例えば、図 1 に関して上述したようなログファイルとして受信および格納することによって、コンテキストデータおよびデータベースオペレーション履歴データを維持する。一実施形態では、データ解析アプリケーション 1 2 5 がデータベースオペレーションを検出すると、ログファイルをデータベースオペレーション履歴モジュール 1 1 2 に送信する。別の実施形態では、データベースオペレーション履歴モジュール 1 1 2 がデータ解析アプリケーション 1 2 5 を継続的に監視し、データベースオペレーションを検出すると、データベースオペレーション履歴データおよび対応するコンテキストデータを受信し、記憶する。

30

40

【 0 0 7 3 】

図 1 に関して上述したように、ガイド付きユーザおよび/またはトレーニングユーザとしてのユーザのステータス、ならびにガイド付きユーザとトレーニングユーザとの組の間の関連付けは、システム管理者、他のユーザによって、または自動的に指定することができる。

【 0 0 7 4 】

ステップ 5 0 5 および 5 1 0 では、ガイド付きユーザに推薦を提供する際に使用するた

50

めに、1つまたは複数の予測モデルが構築され、トレーニングされる。データベースオペレーション推奨モジュール114は、予測モデルを構築する(505)。予測モデルは、演算モデル(OP)、オペランドモデル(OPD)、列演算モデル(OPC)、またはそれらの任意の組合せとすることができる。予測モデルを構築することは、そのデータベースオペレーションがモデルのトレーニングデータとして使用されるトレーニングユーザを決定することを含む。予測モデルを構築することは、モデルクラスを決定することをさらに含む。例えば、予測モデルがOPモデルである場合、クラスはデータベースオペレーションである。予測モデルがOPDモデルである場合、クラスはオペランドである。予測モデルがOPCモデルである場合、クラスは結合および演算、または定義された2テーブル演算である。予測モデルを構築するステップは上記の表5に関して説明したように、可能なモデル特徴を決定するステップをさらに含む。予測モデルを構築するステップは、データベース動作推奨ストア121からモデル方程式を検索するステップをさらに含む。ステップ505の終わりに、モデルはそのトレーニングされていない形式で存在する。図2に関して説明した方程式は各クラスについて組み立てられるが、特徴重みは未知であるか、またはデフォルト値に設定される。この形態では、モデルが決定されたトレーニングユーザに対応する適切なコンテキストデータを用いてトレーニングする準備ができています。

【0075】

モデルトレーニングモジュール210は、決定されたトレーニングユーザからの維持されたデータベースオペレーション履歴データおよびコンテキストデータを使用して、モデルをトレーニングする(510)。モデルトレーニングモジュール210は、プロファイリングデータストア118およびデータベースオペレーション履歴ストア120から、トレーニングユーザに対応するデータベースオペレーション履歴データおよびトレーニングコンテキストデータを検索する。図2に関して上述したように、モデルトレーニングモジュール210は、どのコンテキストデータが特定のデータベースオペレーションまたはオペランドを予測するかを決定する。モデルトレーニングモジュール210は、図2に関して上述したように、各モデル特徴に対する特徴重みを決定する。特徴重みおよび他のパラメータは、データベースオペレーション推奨ストア121に格納され、必要に応じて使用のために取り出され得る。一実施形態では、モデルトレーニングモジュール210が図2に関して上述したように、モデルをトレーニングする前にコンテキストデータを前処理する。一旦、モデルがトレーニングされると、モデルは、特徴のセット(データ解析アプリケーションから受信したコンテキストデータ)に基づいてクラス(オペレーションまたはオペランド)の確率を決定するために使用され得る。

【0076】

ステップ505および510は定期的な間隔で、継続的に、またはどれだけの新しいトレーニングデータが利用可能であるかなどの要因に応じて、行われ得る。ステップ505および510は、データベース動作推奨モジュール114によって生成される各予測モデルに対して繰り返されてもよい。図2に関して上述したように、OPDモデルはOPモデルによって決定されたデータベースオペレーションのためのオペランドを決定するために、OPモデルと共に使用することができる。OPDモデルはOPDモデルによって決定された推奨オペランドが決定された演算に対応するように、OPモデルによって決定されたデータベース演算を入力として取り込むことができる。

【0077】

図5Bは、一実施形態による、データ解析アプリケーションのガイド付きユーザにデータベースオペレーションを推薦するためにトレーニングされた予測モデルを使用する方法を示すフローチャートである。推薦生成モジュール220は、ガイド付きユーザのデータ解析アプリケーション125からアプリケーションコンテキストデータを受信する(550)。一実施形態では、アプリケーション・コンテキスト・データがデータ解析アプリケーション内に表示されたテーブル内で選択された列など、データ解析アプリケーション125との検出された対話に回答して受信される。データ解析アプリケーション125は対話を検出し、コンテキストデータを含むアプリケーションログエントリを作成し、アプリ

ケーションログエントリをデータ解析サーバ104に送信する。一実施形態では、推薦生成モジュール220が、アプリケーションログエントリをプロファイルして、トレーニングされた予測モデルへの入力として使用することができるフォーマットで、コンテキストデータを取り込む。

【0078】

推薦生成モジュール220は、コンテキストデータに基づいて、推薦を生成するために使用する1つ以上のモデルを選択する(555)。例えば、コンテキストデータが、プロジェクトが1つのテーブルを有することを示す場合、推薦生成モジュール220は、OPモデルおよびOPDモデルを使用して、推薦を生成する。コンテキストデータが、プロジェクトが複数のテーブルを有することを示す場合、推薦生成モジュール220は、OPモデル、OPDモデル、およびOPCモデルを使用して、推薦を生成する。図2および図5Aに関して上述したように、OPDモデルは、OPモデルの出力を入力として使用して、OPモデルによって決定された推奨演算のリストに対応するオペランドを決定することができる。

10

【0079】

推薦生成モジュール220は、選択された予測モデルおよび受信されたコンテキストデータを使用して、ガイド付きユーザに推薦するためのデータベースオペレーションおよび/またはオペランドのリストを生成する(560)。様々な実施形態では、生成された推奨のリストがOPモデル、OPDモデル、およびOPCモデル、ならびに他の予測モデルのうちの1つまたは複数によって決定される演算およびオペランドを含む。推薦生成モジュール220は、ステップ555で選択された各モデルを使用して、各モデルクラスに関連する確率を決定する。生成された推奨のリストは、決定された確率に基づいている。例えば、OPモデルまたはOPCモデルが使用される場合、推薦生成モジュール220はモデルによって決定されるような多数の最も確からしいデータベースオペレーションを選択し、案内されたユーザに推薦として提供する。OPDモデルも使用される場合、OPモデルによって決定された選択されたデータベース演算は、選択されたデータベース演算のための最も可能性の高いオペランドの数を決定するためにOPDモデルへの入力として使用される。

20

【0080】

推薦生成モジュール220は、ガイド付きユーザに提示するために、推薦のリストをデータ解析アプリケーション125に送信する(535)。一実施形態では、各推奨データベースオペレーションがデータ解析アプリケーション125のデータベースオペレーションを一意に識別するオペレーション識別子を含む。別の実施形態では、各推奨データベースオペレーションがデータ解析アプリケーション125のユーザに提示するためのデータベースオペレーションのテキスト名または説明をさらに含む。データベースオペレーション、オペレーション識別子、ならびにテキスト名および説明は、データベースオペレーション推奨ストア121に格納され、推薦されたデータベースオペレーションをデータ解析アプリケーション125に送る前に、データベースオペレーション推薦モジュール114によって検索され得る。

30

【0081】

図6は、一実施形態による、選択された列に回答して提供される推奨を備えた、図3の例示的なユーザインタフェースを示す。例示的なユーザインタフェースでは、例えばユーザ入力に回答して、列650が選択される。データ解析アプリケーション125は列選択を検出し、データ解析サーバ104に通知する。データ解析サーバ104は、列650の選択に回答して、データ解析アプリケーション125からコンテキストデータを受信する。一実施形態では、データベースオペレーション推薦モジュール114がユーザのステータスを、ユーザデータストア117からの特定のプロジェクトに対するガイド付きユーザと決定し、コンテキストデータをOPモデル(単一の列が選択されるため)およびOPCモデルに渡す。OPモデルは演算のリストを出力し、OPCモデルは、1つ以上のオペランドを出力する。データベースオペレーション推奨モジュール114は、推奨されるデー

40

50

データベースオペレーション、および適切な場合にはオペランドを決定し、その推奨をデータ解析アプリケーション 125 に送る。図 6 の例では、ユーザが異なる方法でフォーマットされた電話番号を含むように見える列を選択している。したがって、提供される 2 つの推奨は、OP モデルによって決定された電話番号をフォーマットする動作と、OPD モデルによって決定された適用する特定の形式のフォーマットのオペランドとを含む。

【0082】

図 7 は、データ解析アプリケーションにおいて、データ解析サーバから受信した推奨データベースオペレーションおよびオペランドを提示するための方法を示すフローチャートである。データ解析アプリケーション 125 は、データ解析サーバ 104 から推奨データベースオペレーションおよびオペランドを受信する (700)。図 5 に関して上述したように、データベース動作は、データ解析アプリケーション 125 のユーザインタフェースに提示するためのテキスト名または記述を含むことができる。UI モジュール 122 は、データ解析サーバ 104 によって提供されるテキスト名および説明を使用して、推奨データベースオペレーションおよびオペランドに対応するユーザインタフェース要素を生成する (710)。UI モジュール 122 は、データ解析アプリケーション 125 のユーザインタフェースを介して、データ解析アプリケーションのユーザに 1 つ以上の推奨データベースオペレーションを提示する (720)。

【0083】

図 6 に戻ると、提案カード 435 は推奨データベースオペレーションを含む。列 650 は、異なる方法でフォーマットされた電話番号を含む。提案カード 435 上の推奨 660 A ~ C は、セルまたは列内の電話番号をフォーマットすることを含む。推奨 650 A および B は、共通のデータベースオペレーション (電話番号のフォーマット) を有するが、異なるオペランド (電話番号の出力フォーマット) を有する。データ解析アプリケーションのユーザは、推奨 660 A ~ C のうちの 1 つを選択して、データに対して指示されたデータベースオペレーションを実行することができる。

【0084】

追加構成の考慮事項

本明細書で説明するシステムは、クラウドベースのコンピュータ実装を含む、単一のコンピュータまたはコンピュータのネットワークを使用して実装することができる。コンピュータは、好ましくは 1 つまたは複数の高性能 CPU および 1 G またはそれ以上のメインメモリ、ならびに 500 Gb から 2 Tb のコンピュータ可読永続ストレージを含み、Linux またはその変形などのオペレーティングシステムを実行するサーバクラスコンピュータである。本明細書で説明するシステムの動作は、コンピュータストレージにインストールされ、本明細書で説明する機能を実行するために、そのようなサーバのプロセッサによって実行されるハードウェアおよびコンピュータプログラムの組み合わせによって制御することができる。システム 100 はネットワークインターフェースおよびプロトコル、データ入力のための入力デバイス、ならびに表示、印刷、または他のデータの提示のための出力デバイスを含む、本明細書で説明される動作に必要な他のハードウェア要素を含むが、これらは実施形態の関連する詳細を不明瞭にすることを避けるために本明細書では示されない。

【0085】

上記の説明のいくつかの部分は、アルゴリズムのプロセスまたは動作に関して実施形態を説明する。これらのアルゴリズムの説明および表現は、データ処理技術の当業者によって一般的に使用され、彼らの作業内容を他の当業者に効果的に伝える。これらの動作は機能的、計算的、または論理的に説明されているが、プロセッサまたは等価の電気回路、マイクロコードなどによって実行される命令を含むコンピュータプログラムによって実施されるものと理解される。さらに、一般性を損なうことなく、これらの機能的オペレーションの配置をモジュールと呼ぶことも便利である場合があることが判明している。説明された動作およびそれらの関連するモジュールは、ソフトウェア、ファームウェア、ハードウェア、またはそれらの任意の組合せで具現化されてもよい。

【0086】

本明細書で使用されるように、用語「モジュール」は、指定された機能を提供するために利用されるコンピュータプログラムロジックを指す。したがって、モジュールは、ハードウェア、ファームウェア、および/またはソフトウェアで実装することができる。一実施形態では、プログラムモジュールが記憶装置に格納され、メモリにロードされ、プロセッサによって実行される。本明細書で説明される物理的構成要素の実施形態は、本明細書で説明されるもの以外の他のおよび/または異なるモジュールを含むことができる。さらに、他の実施形態では、モジュールに起因する機能が他のモジュールまたは異なるモジュールによって実行することができる。さらに、この説明は、明瞭さおよび便宜のために「モジュール」という用語を省略することがある。

10

【0087】

本発明はまた、本明細書における動作を実行するための装置に関する。この装置は、必要な目的のために特別に構築されてもよく、またはコンピュータによってアクセスされるコンピュータ可読媒体上に格納されたコンピュータプログラムによって選択的に起動または再構成される汎用コンピュータを備えてもよい。そのようなコンピュータプログラムはフロッピー（登録商標）ディスク、光ディスク、CD-ROM、磁気-光ディスク、読み取り専用メモリ（ROM）、ランダムアクセスメモリ（RAM）、EPROM、EEPROM、磁気または光カード、特定用途向け集積回路（ASIC）、または電子命令を記憶するのに適した任意のタイプのコンピュータ可読記憶媒体を含む任意のタイプのディスクなどのコンピュータ可読記憶媒体に記憶することができ、それぞれがコンピュータシステムバスに結合されるが、これらに限定されない。さらに、本明細書で言及するコンピュータは単一のプロセッサを含むことができ、または計算能力を高めるために複数のプロセッサ設計を使用するアーキテクチャとすることができる。

20

【0088】

本明細書で使用される「1つの実施形態」または「一実施形態」への言及は、実施形態に関連して説明された特定の要素、機能、構成、または特徴が少なくとも1つの実施形態に含まれることを手段する。明細書の様々な場所における「一実施形態では」という語句の出現は、必ずしもすべてが同じ実施形態を指すとは限らない。

【0089】

本明細書において用いられるとき、「備える（comprises）」、「備える（comprising）」、「含む（includes）」、「含める（including）」、「有する（has）」、「有する（having）」という用語またはそれらの任意の他の活用形は、非限定的な包含をカバーするものとする。例えば、一連の要素を含むプロセス、方法、物品、または装置は、それらの要素のみに必ずしも限定されず、特に明記されていないかあるいはかかるプロセス、方法、物品、または装置に固有の他の要素を含めてもよい。更に、明確に逆のことを表さない限り、「または」は包括的「or」を指し、排他的「or」を意味しない。例えば、条件AまたはBは、以下のいずれか1つによって満たされる：Aが真であり（または存在する）かつBが偽である（または存在しない）、Aが偽であり（または存在しない）かつBが真である（または存在する）、AおよびBの両方が真である（または存在する）。

30

40

【0090】

さらに、「1つの（a）」または「1つの（an）」の使用は、本明細書の実施形態の要素および構成要素を説明するために使用される。これは、単に便宜上、かつ本開示の一般的な意味を与えるためになされる。本明細書は1つまたは少なくとも1つを含めるように読まれるべきであり、複数でないことを意図することが明白でない限り、単数形は複数形も含める。

【0091】

本開示を読めば、当業者は、識別子空間にわたるエンティティの類似性を決定するためのシステムおよびプロセスのためのさらに追加の代替の構造および機能設計を理解するであろう。したがって、特定の実施形態および用途を図示し、説明したが、本発明は、本明

50

細書に開示される正確な構成および構成要素に限定されず、当業者には明らかな様々な修正、変更、および変形が添付の特許請求の範囲に定義される精神および範囲から逸脱することなく、本明細書に開示される方法および装置の構成、動作、および詳細において行われてもよいことを理解されたい。

【図 1】

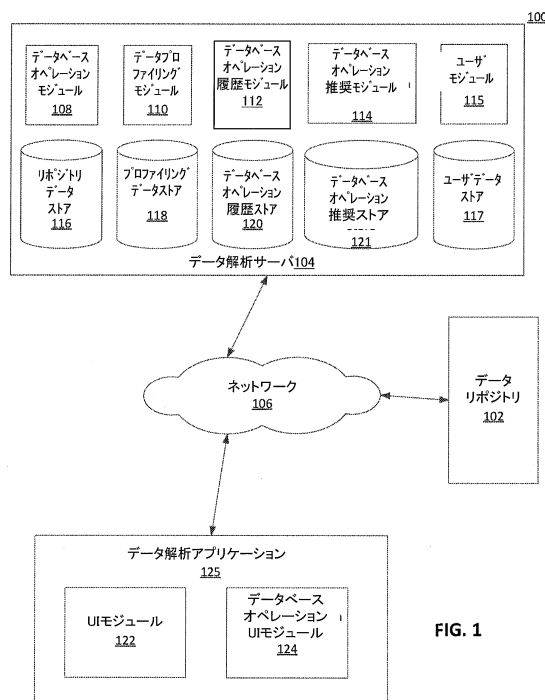


FIG. 1

【図 2】

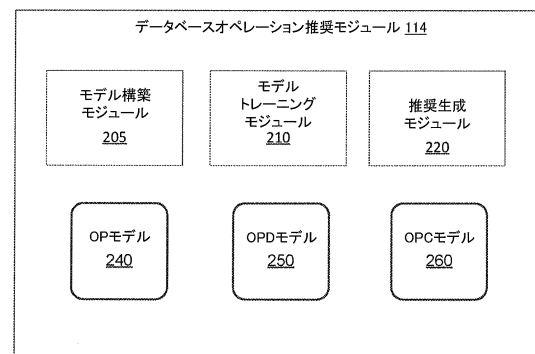


FIG. 2

【 図 3 】

id	type	multis	unique	patterns	domain	macrobase	minibase	operation
L431431	Integer	0	51.71	<NUMBER>	None	9352	12	exp+app
T3607440	String	0	52.81	<WORD>	Handle	-1	-1	valid+domain
E5405405	String	0	100	<WORD><WORD>	email+address	-1	-1	valid+domain
E5405405	String	0	100	<WORD>@<WORD><WORD>	email+address	-1	-1	extra+domain+none
E60319631	String	0	96.79	<WORD>@<WORD><WORD>	email+address	-1	-1	split+on+entity+validation
P40319631	String	0	100	<NUMBER>	None	-1	-1	exp+CONCATENATE
L2152123	String	0	100	<NUMBER>	None	-1	-1	split+text
P4002500	String	0	100	<NUMBER><NUMBER>	None	-1	-1	split+text
E5405405	String	0	100	<WORD><WORD><WORD>	email+address	-1	-1	split+on+entity+validation
L21421236	Disjunctive	0	51.2	<NUMBER><NUMBER><NUMBER>	Date	1,4,2E+12	1,3,9E+12	extra+on+entity+validation
L21421236	Integer	0	0.3	<NUMBER>	None	4	-1	extra+on+entity+validation
O1400130	String	0	0.3	<WORD>	agreement	-1	-1	Agg
S125332	String	0	0.8	<WORD><WORD>	None	-1	-1	Agg
S5905550	String	0	97.79	<NUMBER><NUMBER><NUMBER>	None	-1	-1	split+text

FIG. 3

【 図 4 】

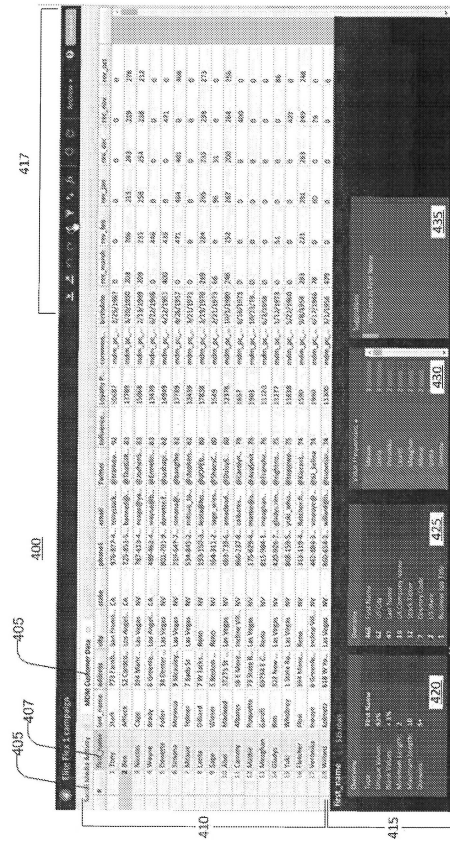


FIG. 4

【 図 5 A 】

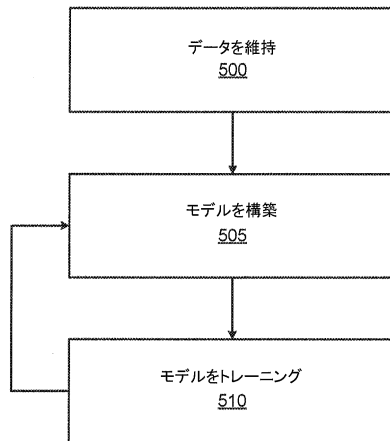


FIG. 5A

【 図 5 B 】

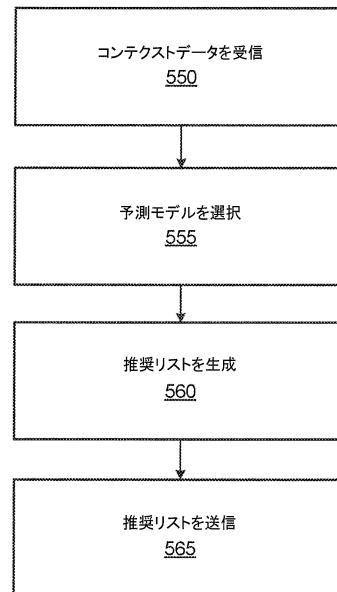


FIG. 5B

【 図 6 】

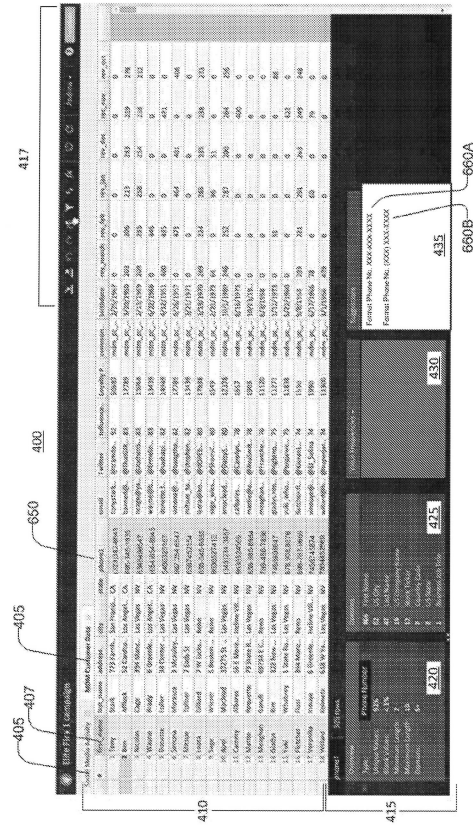


FIG. 6

【圖 7】

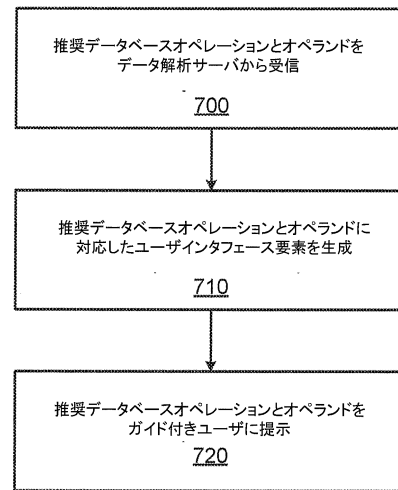


FIG. 7

フロントページの続き

- (72)発明者 デイ、 アトレイー
インド国 560077 バンガロール、 コサヌアー、 バイラシ クロス、 ソリティア レ
ジデンシー B201
- (72)発明者 カルスカル、 サンジェイ
アメリカ合衆国 94063 カリフォルニア州、 レッドウッド シティ、 シーポート ブル
ーバード 2100
- (72)発明者 ダーニング、 ウダヤクマール
アメリカ合衆国 95014 カリフォルニア州、 クパチーノ、 アpartment 645、
バレー グリーン ドライブ 20990

審査官 齊藤 貴孝

- (56)参考文献 米国特許出願公開第2007/0276857(US, A1)
特開2010-033377(JP, A)
中国特許出願公開第101884044(CN, A)
米国特許出願公開第2015/0324346(US, A1)
米国特許出願公開第2016/0092475(US, A1)
濱崎 雅弘, Linked Open DataのためのSPARQLクエリ共有システムの提
案, 一般社団法人 人工知能学会 第29回全国大会論文集CD-ROM [CD-ROM]
2015年度 人工知能学会全国大会(第29回)論文集, 日本, 社団法人人工知能学会, 20
15年 6月 2日, p. 1-4

- (58)調査した分野(Int.Cl., DB名)
G06F 16/00 - 16/958