

【公報種別】特許法第17条の2の規定による補正の掲載
 【部門区分】第6部門第3区分
 【発行日】令和5年6月6日(2023.6.6)

【国際公開番号】WO2020/243499
 【公表番号】特表2022-535792(P2022-535792A)
 【公表日】令和4年8月10日(2022.8.10)
 【年通号数】公開公報(特許)2022-146
 【出願番号】特願2021-571432(P2021-571432)
 【国際特許分類】

10

G 0 6 F 1 6 / 2 6 (2 0 1 9 . 0 1)

【F I】

G 0 6 F 1 6 / 2 6

【手続補正書】

【提出日】令和5年5月29日(2023.5.29)

【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

20

【補正の内容】

【特許請求の範囲】

【請求項1】

1つ又は複数のデータセットに含まれるフィールドのデータ値の意味論的意味を発見するための、データ処理システムによって実行される方法であって、前記方法は、

前記1つ又は複数のデータセットに含まれるフィールドを識別することであって、前記フィールドは、識別子に関連付けられる、識別することと、

前記フィールドの1つ又は複数のデータ値を識別することであって、少なくとも前記データ値のうちの第1のデータ値は、前記データ値のうちの第2のデータ値と異なる、識別することと、

30

前記フィールドの前記1つ又は複数のデータ値に基づいて、前記フィールドの前記1つ又は複数のデータ値の1つ又は複数の属性を特定することと、

複数のテストにアクセスすることであって、テストは、1つ又は複数の所定の属性に対して提案された意味論的意味を指定する、アクセスすることと、

意味論的意味は、所定のフィールドにどのような種類のデータ値が含まれるかを示し、

前記フィールドに対して前記特定された1つ又は複数の属性に少なくとも前記複数のテストを適用することに基づいて、前記フィールドの前記1つ又は複数のデータ値の1つ又は複数の提案された意味論的意味を指定することと、

前記フィールドに対して前記1つ又は複数の提案された意味論的意味の間の類似性を特定することと、

40

前記1つ又は複数の提案された意味論的意味の間の前記類似性に基づいて、前記提案された意味論的意味の1つを、前記フィールドの前記1つ又は複数のデータ値の前記意味論的意味を識別するものとして識別することと、

データストア内において、前記フィールドの前記識別子を、前記フィールドの前記1つ又は複数のデータ値の前記識別された意味論的意味と共に保存することとを含む、方法。

【請求項2】

前記フィールドのデータ値のフォーマットを特定することをさらに含み、前記1つ又は複数の属性は前記フォーマットを示す、請求項1に記載の方法。

【請求項3】

50

前記フィールドに含まれる前記データ値を表す統計値を特定することをさらに含み、前記1つ又は複数の属性は統計値を示す、請求項1に記載の方法。

【請求項4】

前記統計値は、前記フィールドの前記1つ又は複数のデータ値の最小長さ、前記フィールドの前記1つ又は複数のデータ値の最大長さ、前記フィールドの最も一般的なデータ値、前記フィールドの最も一般的でないデータ値、前記フィールドの最大データ値及び前記フィールドの最小データ値の少なくとも1つを含む、請求項3に記載の方法。

【請求項5】

前記複数のテストを適用することは、
前記フィールドが、前記1つ又は複数のデータセットのうちのデータセットのためのプライマリキーを含むことを特定することと、
前記プライマリキーに関連する前記複数のテストのうち1つのテストを選択することとを含む、請求項1に記載の方法。 10

【請求項6】

前記複数のテストを適用することは、前記フィールドのデータ値の、用語集内の用語とのメタデータ比較を行うことを含む、請求項1に記載の方法。

【請求項7】

前記複数のテストを適用することは、
前記フィールドの前記データ値の前記1つ又は複数の属性から、前記フィールドの前記1つ又は複数のデータ値によって表されるパターンを特定することと、
前記パターンにマッピングされる特定の意味論的意味を特定することと、
前記フィールドを前記特定の意味論的意味でラベル付けすることとを含む、請求項1に記載の方法。 20

【請求項8】

前記複数のテストを適用することは、
データ集合を表す値のリストを検索することと、
前記フィールドの前記1つ又は複数のデータ値を前記値のリストと比較することと、
前記比較することに応答して、前記データ値の閾値数が値の前記リストの前記値と一致することを特定することと、
前記特定することに応答して、前記フィールドを、前記データ集合を指定する特定の意味論的意味でラベル付けすることとを含む、請求項1に記載の方法。 30

【請求項9】

前記複数のテストを適用することは、
前記フィールドのための少なくとも2つの意味論的意味を生成することと、
前記少なくとも2つの意味論的意味が相互に排他的であるか又は包含的であるかを特定することとを含む、請求項1に記載の方法。

【請求項10】

前記複数のテストを適用することに応答して、前記フィールドと、前記1つ又は複数のデータセットの別のフィールドとの間の関係を特定することをさらに含む、請求項1に記載の方法。 40

【請求項11】

前記関係は、前記フィールドの第一のデータ値が、他のフィールドに保存された第二のデータ値を決定するという表示、前記第一のデータ値が前記第二のデータ値と相関するという表示又は前記第一のデータ値が前記第二のデータ値と同一であるという表示の1つを含む、請求項10に記載の方法。

【請求項12】

前記複数のテストは、少なくとも1つの重み値にそれぞれ関連付けられ、前記方法は、少なくとも1つのテストに関連付けられた重み値を更新することと、 50

前記更新された重み値を使用して、前記フィールドの前記データ値の前記1つ又は複数の属性に基づいて前記テストを再び適用することとをさらに含む、請求項1に記載の方法。

【請求項13】

機械学習プロセスを使用して前記複数のテストを訓練することをさらに含む、請求項1に記載の方法。

【請求項14】

データ品質ルール環境から、前記提案された意味論的意味に割り当てられる1つ又は複数のデータ品質ルールを検索することと、

前記1つ又は複数のデータ品質ルールのうちのデータ品質ルールを前記フィールドに割り当てることと

をさらに含む、請求項1に記載の方法。

【請求項15】

スコア値を各提案された意味論的意味に適用することと、

前記提案された意味論的意味の各意味論的意味について、前記意味論的意味に関連付けられた前記スコア値を結合することと、

各提案された意味論的意味に関連付けられた前記スコア値に従って前記提案された意味論的意味をランク付けすることと

をさらに含む、請求項1に記載の方法。

【請求項16】

前記複数のテストから前記提案された意味論的意味のバリデーションを受け取ることと、

前記バリデーションを受け取ることに応答して、前記複数のテストを重み付けすることと

をさらに含む、請求項1に記載の方法。

【請求項17】

前記データストアは、データ辞書を含む、請求項1に記載の方法。

【請求項18】

前記提案された意味論的意味をデータ品質ルール環境に出力することをさらに含む、請求項1に記載の方法。

【請求項19】

前記提案された意味論的意味の前記識別された1つに基づいて、前記データ品質ルール環境からのデータ品質ルールを使用して前記フィールドのデータを処理することに関するエラーの数を、前記提案された意味論的意味の前記識別された1つを使用することなく前記フィールドの前記データを処理することに関するエラーの別の数に対して減少させることをさらに含む、請求項18に記載の方法。

【請求項20】

特定の識別子を含む前記1つ又は複数のデータセットを処理するためのリクエストを受信することと、

前記リクエストに応答して、前記データストアから、前記特定の識別子と、前記特定の識別子に関連付けられた特定の意味論的意味とにアクセスすることと、

前記特定の意味論的意味に基づいて、前記1つ又は複数のデータセットの少なくとも一部を処理するための1つ又は複数のデータ処理ルールを特定することと、

少なくとも前記1つ又は複数のデータセットの前記一部を、前記選択された1つ又は複数のデータ処理ルールに従って処理することと

をさらに含む、請求項1に記載の方法。

【請求項21】

1つ又は複数のデータセットに含まれるフィールドの意味論的意味を発見するためのデータ処理システムであって、前記データ処理システムは、

命令を保存するデータストレージと、

10

20

30

40

50

前記データストレージによって保存された前記命令を実行して、

1つ又は複数のデータセットに含まれるフィールドを識別することであって、前記フィールドは、識別子に関連付けられている、識別することと、

前記フィールドの1つ又は複数のデータ値を識別することであって、少なくとも前記データ値のうちの第1のデータ値は、前記データ値のうちの第2のデータ値と異なる、識別することと、

前記フィールドの前記1つ又は複数のデータ値に基づいて、前記フィールドの前記データ値の1つ又は複数の属性を特定することと、

複数のテストにアクセスすることであって、テストは、1つ又は複数の所定の属性に対して提案された意味論的意味を指定する、アクセスすることと、

10

意味論的意味は、所定のフィールドにどのような種類のデータ値が含まれるかを示し、

前記フィールドに対して前記特定された1つ又は複数の属性に少なくとも前記複数のテストを適用することに基づいて、前記フィールドの前記1つ又は複数のデータ値の1つ又は複数の提案された意味論的意味を指定することと、

前記フィールドに対して前記1つ又は複数の提案された意味論的意味の間の類似性を特定することと、

前記1つ又は複数の提案された意味論的意味の間の前記類似性に基づいて、前記提案された意味論的意味の1つを、前記フィールドの前記1つ又は複数のデータ値の前記意味論的意味を識別するものとして識別することと、

データストア内において、前記フィールドの前記識別子を、前記フィールドの前記1つ又は複数のデータ値の前記識別された意味論的意味と共に保存することと

20

を含む動作を行うように構成された少なくとも1つのプロセッサとを含むデータ処理システム。

【請求項22】

前記動作は、

特定の識別子を含む前記1つ又は複数のデータセットを処理するためのリクエストを受信することと、

前記リクエストに応答して、前記データストアから、前記特定の識別子と、前記特定の識別子に関連付けられた特定の意味論的意味とにアクセスすることと、

前記特定の意味論的意味に基づいて、前記1つ又は複数のデータセットの少なくとも一部を処理するための1つ又は複数のデータ処理ルールを特定することと、

30

少なくとも前記1つ又は複数のデータセットの前記一部を、前記選択された1つ又は複数のデータ処理ルールに従って処理することと

をさらに含む請求項21に記載のデータ処理システム。

【請求項23】

前記動作は、前記フィールドのデータ値のフォーマットを特定することをさらに含み、前記1つ又は複数の属性は前記フォーマットを示す、請求項21に記載のデータ処理システム。

【請求項24】

前記動作は、前記フィールドに含まれる前記データ値を表す統計値を特定することをさらに含み、前記1つ又は複数の属性は統計値を示す、請求項21に記載のデータ処理システム。

40

【請求項25】

前記統計値は、前記フィールドの前記1つ又は複数のデータ値の最小長さ、前記フィールドの前記1つ又は複数のデータ値の最大長さ、前記フィールドの最も一般的なデータ値、前記フィールドの最も一般的でないデータ値、前記フィールドの最大データ値及び前記フィールドの最小データ値の少なくとも1つを含む、請求項24に記載のデータ処理システム。

【請求項26】

前記複数のテストを適用することは、

50

前記フィールドが、前記1つ又は複数のデータセットのうちの前記データセットのためのプライマリーを含むことを特定することと、

前記プライマリーに関連する前記複数のテストのうち1つのテストを選択することとを含む、請求項21に記載のデータ処理システム。

【請求項27】

前記複数のテストを適用することは、前記フィールドのデータ値の、用語集内の用語とのメタデータ比較を行うことを含む、請求項21に記載のデータ処理システム。

【請求項28】

前記複数のテストを適用することは、

前記フィールドの前記データ値の前記1つ又は複数の属性から、前記フィールドの前記1つ又は複数のデータ値によって表されるパターンを特定することと、

前記パターンにマッピングされる特定の意味論的意味を特定することと、

前記フィールドを前記特定の意味論的意味でラベル付けすることと

を含む、請求項21に記載のデータ処理システム。

【請求項29】

前記複数のテストを適用することは、

データ集合を表す値のリストを検索することと、

前記フィールドの前記1つ又は複数のデータ値を前記値のリストと比較することと、

前記比較することに応答して、前記データ値の閾値数が値の前記リストの前記値と一致することを特定することと、

前記特定することに応答して、前記フィールドを、前記データ集合を指定する特定の意味論的意味でラベル付けすることと

を含む、請求項21に記載のデータ処理システム。

【請求項30】

1つ又は複数のデータセットに含まれるフィールドの意味論的意味を発見するための命令を記憶する1つ又は複数の非一時的コンピュータ可読媒体であって、前記命令は、

1つ又は複数のデータセットに含まれるフィールドを識別することであって、前記フィールドは、識別子に関連付けられている、識別することと、

前記フィールドの1つ又は複数のデータ値を識別することであって、少なくとも前記データ値のうちの前記第1のデータ値は、前記データ値のうちの前記第2のデータ値と異なる、識別することと、

前記フィールドの前記1つ又は複数のデータ値に基づいて、前記フィールドの前記データ値の1つ又は複数の属性を特定することと、

複数のテストにアクセスすることであって、テストは、1つ又は複数の所定の属性に対して提案された意味論的意味を指定する、アクセスすることと、

意味論的意味は、所定のフィールドにどのような種類のデータ値が含まれるかを示し、

前記フィールドに対して前記特定された1つ又は複数の属性に少なくとも前記複数のテストを適用することに基づいて、前記フィールドの前記1つ又は複数のデータ値の1つ又は複数の提案された意味論的意味を指定することと、

前記フィールドに対して前記1つ又は複数の提案された意味論的意味の間の類似性を特定することと、

前記1つ又は複数の提案された意味論的意味の間の前記類似性に基づいて、前記提案された意味論的意味の1つを、前記フィールドの前記1つ又は複数のデータ値の前記意味論的意味を識別するものとして識別することと、

データストア内において、前記フィールドの前記識別子を、前記フィールドの前記1つ又は複数のデータ値の前記識別された意味論的意味と共に保存することと

を含む動作を行うように構成された1つ又は複数のプロセッサによって実行可能である、1つ又は複数の非一時的コンピュータ可読媒体。

【請求項31】

前記動作は、

10

20

30

40

50

特定の識別子を含む前記 1 つ又は複数のデータセットを処理するためのリクエストを受信することと、

前記リクエストに応答して、前記データストアから、前記特定の識別子と、前記特定の識別子に関連付けられた特定の意味論的意味とにアクセスすることと、

前記特定の意味論的意味に基づいて、前記 1 つ又は複数のデータセットの少なくとも一部を処理するための 1 つ又は複数のデータ処理ルールを特定することと、

少なくとも前記 1 つ又は複数のデータセットの前記一部を、前記選択された 1 つ又は複数のデータ処理ルールに従って処理することと

をさらに含む、請求項 30 に記載の 1 つ又は複数の非一時的コンピュータ可読媒体。

10

20

30

40

50