

(19) 日本国特許庁(JP)

(12) 特許公報(B2)

(11) 特許番号

特許第4687089号
(P4687089)

(45) 発行日 平成23年5月25日(2011.5.25)

(24) 登録日 平成23年2月25日(2011.2.25)

(51) Int.Cl.	F I
G06F 17/30 (2006.01)	G06F 17/30 180D
	G06F 17/30 350C
	G06F 17/30 320D

請求項の数 14 (全 27 頁)

(21) 出願番号	特願2004-355789 (P2004-355789)	(73) 特許権者	000004237
(22) 出願日	平成16年12月8日(2004.12.8)		日本電気株式会社
(65) 公開番号	特開2006-163941 (P2006-163941A)		東京都港区芝五丁目7番1号
(43) 公開日	平成18年6月22日(2006.6.22)	(74) 代理人	100103090
審査請求日	平成19年11月12日(2007.11.12)		弁理士 岩壁 冬樹
		(74) 代理人	100124501
			弁理士 塩川 誠人
		(72) 発明者	久寿居 大
			東京都港区芝五丁目7番1号 日本電気株式会社内
		(72) 発明者	立石 健二
			東京都港区芝五丁目7番1号 日本電気株式会社内

最終頁に続く

(54) 【発明の名称】 重複レコード検出システム、および重複レコード検出プログラム

(57) 【特許請求の範囲】

【請求項1】

語の変換に用いられる辞書であって、当該語に対応する代表的な語である代表語を対応付けた辞書である代表語辞書と、相互に省略可能な前記代表語である省略可能語をグループ化した辞書である省略語辞書とを記憶する変換語記憶部と、

複数の情報からなる複数のレコードを保持するデータベースの各レコード間の表記の類似度を計算する類似度計算部と、

前記類似度計算部が計算した前記類似度が所定の値以上であるレコードの組み合わせである重複レコード候補を抽出する重複候補抽出部とを備え、

前記類似度計算部は、前記各レコードに含まれる語のうち、前記代表語辞書に含まれる語を対応する代表語に変換し、当該代表語に隣接する位置に前記省略語辞書において当該代表語と同一のグループに含まれる代表語を追加し、代表語が追加された各レコード間の表記の類似度を計算する

ことを特徴とする重複レコード検出システム。

【請求項2】

変換語記憶部は、同義語を代表語として記憶し、

類似度計算部は、データベースに登録されている各レコードに含まれる語を、対応する同義語に変換してレコード間の表記の類似度を計算する

請求項1記載の重複レコード検出システム。

【請求項3】

10

20

データベースのレコードを構成し、前記データベースのレコードに登録されている情報が区切られる単位であるフィールドの情報を入力するデータベース情報入力部を含み、

類似度計算部は、各レコード間の表記の類似度を前記フィールドごとに算出し、入力されたフィールドの情報に基づいて、フィールドごとに算出された類似度から、各レコード間の表記の類似度を計算する

請求項 1 または請求項 2 記載の重複レコード検出システム。

【請求項 4】

重複候補抽出部が抽出した重複レコード候補が互いに同一の内容の情報のレコードの組み合わせである重複レコードであるか否かを類似度に応じて規定したルールである重複判定ルールを記憶する重複判定ルール記憶部と、

前記重複判定ルール記憶部が記憶している前記重複判定ルールに規定された類似度と、各レコード間の類似度または各レコード間のフィールドごとの類似度との関係にもとづいて、前記重複レコード候補が前記重複レコードであるか否かを判定する重複レコード判定部とを含む

請求項 1 から請求項 3 のうちいずれか 1 項に記載の重複レコード検出システム。

【請求項 5】

重複レコード判定部が重複レコードであると判定したレコード間で、異なる部分から導出される語の組を代表語候補として抽出し、抽出した代表語候補を代表語辞書に含めて変換語記憶部に記憶させる代表語候補抽出部を含む

請求項 4 記載の重複レコード検出システム。

【請求項 6】

代表語候補抽出部は、重複レコードと判定された 2 つのレコードのうち、一のレコードの文字列が他のレコードの文字列に含まれる場合、2 つのレコードで異なる部分の文字列と、共通する部分の文字列との組を、省略可能語候補として抽出し、抽出した省略可能語候補を省略語辞書に含めて変換語記憶部に記憶させる

請求項 5 記載の重複レコード検出システム。

【請求項 7】

代表語候補抽出部は、データベース内のレコードのうちいずれかのレコードにおいて、抽出された代表語候補に含まれる語を全て含むレコードが存在する場合、当該代表語候補を、省略可能語候補とし、当該省略可能語候補を省略語辞書に含めて変換語記憶部に記憶させる

請求項 5 記載の重複レコード検出システム。

【請求項 8】

代表語候補抽出部は、抽出した代表語候補のうち、当該代表語候補に含まれる語を組み合わせた文字列が、他の代表語候補に含まれる語と一致する場合、当該代表語候補を、変換語記憶部に記憶させる対象の代表語候補から除外する

請求項 7 記載の重複レコード検出システム。

【請求項 9】

代表語候補抽出部は、抽出された省略可能語候補がいずれかの代表語候補に含まれる組の語をいずれも含んでいる場合、当該省略可能語候補を、変換語記憶部に記憶させる対象の代表語から除外する

請求項 8 記載の重複レコード検出システム。

【請求項 10】

重複レコード候補を出力し、使用者が、重複レコード候補が重複レコードであるか否かの判定を入力する入出力部を含み、

重複レコード判定部は、重複候補抽出部が抽出した重複レコード候補のうち、重複判定ルールにより重複レコードでないと判定された重複レコード候補を、当該重複レコード候補に含まれる語の数の最も多い重複レコード候補から順に前記入出力部に出力する

請求項 4 から請求項 9 のうちいずれか 1 項記載の重複レコード検出システム。

【請求項 11】

10

20

30

40

50

重複レコード判定部は、重複判定ルールにより重複レコードでないと判定された重複レコード候補のうち、重複すると判定されるレコードの組合せの数が多い順に、前記重複レコード候補を入出力部に出力する

請求項 1 0 記載の重複レコード検出システム。

【請求項 1 2】

重複レコードであると重複レコード判定部、または入出力部を介して使用者に判定されたレコードの組み合わせを記憶する重複レコードデータベースと、

前記重複レコードデータベースが記憶しているレコードの組み合わせを構成するレコードのうち、一のレコード以外のレコードをデータベースから削除する重複レコード削除部とを含む

10

請求項 1 0 または請求項 1 1 記載の重複レコード検出システム。

【請求項 1 3】

データベースに登録すべく使用者が入力した情報の語を、変換語記憶部が記憶している語に変換して、変換した語、または入力された前記情報の語と合致する語からなる情報を含むレコードをデータベースから抽出するデータベース登録部と、

前記データベース登録部が抽出したレコードを表示する表示部とを含む

請求項 1 から請求項 1 2 のうちいずれか 1 項記載の重複レコード検出システム。

【請求項 1 4】

語の変換に用いられる辞書であって、当該語に対応する代表的な語である代表語を対応付けた辞書である代表語辞書と、相互に省略可能な前記代表語である省略可能語をグループ化した辞書である省略語辞書とを記憶する変換語記憶部を備えたコンピュータに適用される重複レコード検出プログラムであって、

20

前記コンピュータに、

複数の情報からなる複数のレコードを保持するデータベースの各レコード間の表記の類似度を計算する類似度計算処理と、

前記類似度計算処理で計算された前記類似度が、所定の値以上であるレコードの組み合わせである重複レコード候補を抽出する重複候補抽出処理とを実行させ、

前記類似度計算処理で、前記各レコードに含まれる語のうち、前記代表語辞書に含まれる語を対応する代表語に変換させ、当該代表語に隣接する位置に前記省略語辞書において当該代表語と同一のグループに含まれる代表語を追加させ、代表語が追加された各レコード間の表記の類似度を計算させる

30

ための重複レコード検出プログラム。

【発明の詳細な説明】

【技術分野】

【0001】

本発明は、店舗等の情報が登録されたデータベースに重複して登録されている情報を検出する重複レコード検出システム、および重複レコード検出プログラムに関する。

【背景技術】

【0002】

店舗や、人物、書物等の情報によって構成されるデータベースに、重複する情報が登録されている場合がある。

40

【0003】

具体的には、例えば、同じ情報が異なる書式でデータベースに登録されていたり、同義であるが異なる語によってデータベースに登録されていたりする。同じ情報が重複してデータベースに登録されていると、データベースの容量が大きくなってしまったりするという問題がある。

【0004】

データベースの容量を削減するために、特許文献 1 には、多数の人物の情報が登録されているデータベースから、同一人物の情報の重複登録を検出するシステムが記載されている。

50

【 0 0 0 5 】

また、特許文献 2 には、書物の情報が登録されているデータベースから一の書物を検索対象として抽出する際に、異なる書式によって同一の書物が複数重複してデータベースに登録されていても、検索対象となる書物をすべて抽出する装置が記載されている。

【 0 0 0 6 】

【特許文献 1】特開平 1 1 - 1 8 4 8 8 4 号公報 (段落 0 0 1 7 ~ 0 0 4 9、図 1)

【特許文献 2】特開 2 0 0 4 - 2 9 9 6 9 号公報 (段落 0 0 2 2 ~ 0 0 7 5、図 2)

【発明の開示】

【発明が解決しようとする課題】

【 0 0 0 7 】

特許文献 1 に記載されているシステムは、例えば、カナ氏名、漢字氏名、カナ住所、漢字住所、生年月日などの書式を、統一した表記法による書式に正規化することによって、同一人物の情報の重複登録を検出する。

【 0 0 0 8 】

しかし、特許文献 1 に記載されているシステムは、同義であるが異なる語によってデータベースに登録されている情報の重複登録を検出することができないという問題がある。具体的には、例えば、同一人物の生年月日が西暦と和暦とで重複して登録されていると、重複登録を検出することができない。また、同一人物の住所の情報が、「東京都千代田区・・」という表記と、「都内千代田区・・」という表記とで重複して登録されていると、重複登録を検出することができない。

【 0 0 0 9 】

また、特許文献 2 に記載されている装置は、入力された検索対象の書物の情報と、データベースに登録されている書物の情報との類似度を算出して類似度の高い書物を検索結果として抽出するが、同義であるが異なる語によってデータベースに登録されている情報の類似度は低くなるため、そのような書物は抽出されにくいという問題がある。

【 0 0 1 0 】

具体的には、例えば、検索対象の書物の名称として「にほん」と入力された場合、「にっぽん」という名称の書物の類似度は低く算出されるため、「にっぽん」という名称の書物は抽出されにくくなってしまう。

【 0 0 1 1 】

そこで、本発明は、同義語や省略可能語による表記の差異があっても、重複する情報をデータベースから検出する重複レコード検出システム、および重複レコード検出プログラムを提供することを目的とする。

【課題を解決するための手段】

【 0 0 1 2 】

本発明による重複レコード検出システムは、語の変換に用いられる辞書であって、その語に対応する代表的な語である代表語を対応付けた辞書である代表語辞書と、相互に省略可能な代表語である省略可能語をグループ化した辞書である省略語辞書とを記憶する変換語記憶部と、複数の情報からなる複数のレコードを保持するデータベースの各レコード間の表記の類似度を計算する類似度計算部と、類似度計算部が計算した類似度が所定の値以上であるレコードの組み合わせである重複レコード候補を抽出する重複候補抽出部とを備え、類似度計算部が、各レコードに含まれる語のうち、代表語辞書に含まれる語に対応する代表語に変換し、当該代表語に隣接する位置に省略語辞書において当該代表語と同一のグループに含まれる代表語を追加し、代表語が追加された各レコード間の表記の類似度を計算することを特徴とする。

【 0 0 1 3 】

変換語記憶部は、同義語を代表語として記憶してもよく、類似度計算部は、データベースに登録されている各レコードに含まれる語を、対応する同義語に変換してレコード間の表記の類似度を計算してもよい。

【 0 0 1 5 】

10

20

30

40

50

データベースのレコードを構成し、データベースのレコードに登録されている情報が区切られる単位であるフィールドの情報を入力するデータベース情報入力部を含んでもよく、類似度計算部は、各レコード間の表記の類似度をフィールドごとに算出し、入力されたフィールドの情報に基づいて、フィールドごとに算出された類似度から、各レコード間の表記の類似度を計算してもよい。そのような構成によれば、フィールド間の類似度を用いて、レコード間の類似度を計算することができる。

【0016】

重複候補抽出部が抽出した重複レコード候補が、互いに同一の内容の情報のレコードの組み合わせである重複レコードであるか否かを類似度に応じて規定したルールである重複判定ルールを記憶する重複判定ルール記憶部と、重複判定ルール記憶部が記憶している重複判定ルールに規定された類似度と、各レコード間の類似度または各レコード間のフィールドごとの類似度との関係にもとづいて、重複レコード候補が重複レコードであるか否かを判定する重複レコード判定部とを含んでもよい。そのような構成によれば、重複レコード候補が重複レコードであるか否かを、自動的に判定することができる。

10

【0017】

重複レコード判定部が重複レコードであると判定したレコード間で、異なる部分から導出される語の組を代表語候補として抽出し、抽出した代表語候補を代表語辞書に含めて変換語記憶部に記憶させる代表語候補抽出部を含んでもよい。そのような構成によれば、重複レコードから、変換語候補を抽出することができる。

【0018】

代表語候補抽出部は、重複レコードと判定された2つのレコードのうち、一のレコードの文字列が他のレコードの文字列に含まれる場合、2つのレコードで異なる部分の文字列と、共通する部分の文字列との組を、省略可能語候補として抽出し、抽出した省略可能語候補を省略語辞書に含めて変換語記憶部に記憶させてもよい。

20

【0019】

代表語候補抽出部は、データベース内のレコードのうちいずれかのレコードにおいて、抽出された代表語候補に含まれる語を全て含むレコードが存在する場合、当該代表語候補を、省略可能語候補とし、当該省略可能語候補を省略語辞書に含めて変換語記憶部に記憶させてもよい。また、代表語候補抽出部は、抽出した代表語候補のうち、当該代表語候補に含まれる語を組み合わせた文字列が、他の代表語候補に含まれる語と一致する場合、当該代表語候補を、変換語記憶部に記憶させる対象の代表語候補から除外してもよい。また、代表語候補抽出部は、抽出された省略可能語候補がいずれかの代表語候補に含まれる組の語をいずれも含んでいる場合、当該省略可能語候補を、変換語記憶部に記憶させる対象の代表語から除外してもよい。

30

【0021】

重複レコード候補を出力し、使用者が、重複レコード候補が重複レコードであるか否かの判定を入力する入出力部を含んでもよく、重複レコード判定部は、重複候補抽出部が抽出した重複レコード候補のうち、重複判定ルールにより重複レコードでないと判定された重複レコード候補を、当該重複レコード候補に含まれる語の数の最も多い重複レコード候補から順に前記入出力部に出力してもよい。そのような構成によれば、使用者が、入出力部を介して判定を入力する回数を減らすことができる。

40

【0022】

重複レコード判定部は、重複判定ルールにより重複レコードでないと判定された重複レコード候補のうち、重複すると判定されるレコードの組合せの数が多い順に、重複レコード候補を入出力部に出力してもよい。

【0024】

重複レコードであると重複レコード判定部、または入出力部を介して使用者に判定されたレコードの組み合わせを記憶する重複レコードデータベースと、重複レコードデータベースが記憶しているレコードの組み合わせを構成するレコードのうち、一のレコード以外のレコードをデータベースから削除する重複レコード削除部とを含んでもよい。そのよう

50

な構成によれば、重複レコードをデータベースから削除することができる。

【0025】

データベースに登録すべく使用者が入力した情報の語を、変換語記憶部が記憶している語に変換して、変換した語、または入力された情報の語と合致する語からなる情報を含むレコードをデータベースから抽出するデータベース登録部と、データベース登録部が抽出したレコードを表示する表示部とを含んでもよい。そのような構成によれば、重複する情報のデータベースへの登録を防ぐことができる。

【0026】

本発明による重複レコード検出プログラムは、語の変換に用いられる辞書であって、当該語に対応する代表的な語である代表語を対応付けた辞書である代表語辞書と、相互に省略可能な前記代表語である省略可能語をグループ化した辞書である省略語辞書とを記憶する変換語記憶部を備えたコンピュータに適用される重複レコード検出プログラムであって、コンピュータに、複数の情報からなる複数のレコードを保持するデータベースの各レコード間の表記の類似度を計算する類似度計算処理と、類似度計算処理で計算した類似度が、所定の値以上であるレコードの組み合わせである重複レコード候補を抽出する重複候補抽出処理とを実行させ、類似度計算処理で、各レコードに含まれる語のうち、代表語辞書に含まれる語を対応する代表語に変換させ、その代表語に隣接する位置に省略語辞書においてその代表語と同一のグループに含まれる代表語を追加させ、代表語が追加された各レコード間の表記の類似度を計算させることを特徴とする。

【発明の効果】

【0027】

本発明によれば、同義語や省略可能語による表記の差異があっても、重複する情報をデータベースから検出することができる。

【発明を実施するための最良の形態】

【0028】

実施の形態1.

本発明の第1の実施の形態について、図面を参照して説明する。図1は、本発明の第1の実施の形態の一構成例を説明するブロック図である。

【0029】

本発明の第1の実施の形態による重複レコード検出システム20は、同義語が登録されている同義語辞書と、省略可能な語が登録されている省略可能語辞書とによって構成される変換語辞書(変換語記憶部)5、変換語辞書5を用いて、店舗の情報が登録されているデータベース2に登録されている複数のレコード間の類似度を計算する類似度計算部3、類似度計算部3が計算した類似度が所定の閾値以上であった情報を抽出する重複候補抽出部6、および各部の動作を制御するプログラムを記憶する記憶部1を含む。

【0030】

重複レコード検出システム20は、プログラムによって処理を実行するサーバ等のコンピュータによって実現される。なお、重複レコード検出システム20は、外部の記憶媒体が記憶しているプログラムに従って処理を実行してもよい。また、変換語辞書5は、予め同義語および省略可能語が登録されているものとする。

【0031】

図2は、データベース2に登録されている情報の例を示す説明図である。データベース2には、例えば、店舗の名称や住所、電話番号が登録されている。なお、データベース2には、各店舗の情報が、レコードに区切られて登録され、店舗の各情報は、登録されているレコードの各情報の属性に応じたフィールドに区切られて登録されているものとする。具体的には、図2に示す例では、レコードIDが「001」のレコードには「エヌイーシー奈良支店」の各情報が登録され、レコードIDが「002」のレコードには「日電奈良支店」の各情報が登録され、レコードIDが「003」のレコードには「NEC奈良支店」の各情報が登録されている。

【0032】

10

20

30

40

50

また、データベース2の各レコードの名称のフィールドには、「エヌイーシー奈良支店」、「日電奈良支店」、および「NEC奈良支店」が登録されており、住所のフィールドには、「1の1」、「1-1」、および「1-1」が登録されており、電話番号のフィールドには、「000-111-1234」、「000-111-1235」、および「000-111-1234」が登録されているものとする。

【0033】

なお、データベース2における各レコードのフィールドの数等の情報を入力し、入力されたデータベース2の情報を類似度計算部3に出力するデータベース情報入力部4を含んでもよい。データベース情報入力部4には、使用者がキーボード等の入力手段を用いてデータベース2の情報を入力してもよい。また、データベース情報入力部4は、記憶部1や外部の記憶媒体が記憶しているデータベース2の情報を読み込んでもよい。

10

【0034】

データベース情報入力部4には、例えば、どのフィールドは何を表しているのか（例えば、各レコードの先頭のフィールドはIDである等）、どのフィールドとどのフィールドとを結合して1つのフィールドとして扱う（例えば、住所が「都道府県」、「市町村」、および「番地とビル名」に分かれている各フィールドを1つのフィールドとして扱う等）のか、どのフィールドの類似度計算に変換語辞書5を用いるのか、およびレコード間の類似度を算出する際の各フィールドの重み（名称フィールド、住所フィールド、および電話番号フィールドの重みの比を、1:1:1とする）等の情報を入力する。

【0035】

20

図3は、変換語辞書5に登録されている情報の例を示す説明図である。変換語辞書5を構成する同義語辞書には、語と、その語の同義語のうち代表的な語である代表語とが対応づけられて登録されている。また、変換語辞書5を構成する省略可能語辞書には、代表語のうち、相互に省略可能な代表語に同じグループIDが付されて登録されている。

【0036】

図3の例によれば、「エヌイーシー」と「日本電気」との代表語は「NEC」であり、「日本電気株式会社」の代表語は「日電」である。また、「NEC」と「日電」とは相互に省略可能な省略可能語である。

【0037】

なお、同義語辞書において、同義語の欄の語は、代表語として用いられることはないものとする。また、省略可能語辞書において、省略可能語として登録されている語は、代表語であってもよいが、同義語ではないものとする。

30

【0038】

類似度計算部3は、例えば、形態素解析等の方法を用いて、データベース2に登録されている各情報を語の単位に分解する。なお、情報を語の単位に分解する他の方法として、例えば、スペースの前後で語の単位に分解したり、文字種が切り替わる位置（例えば、カタカナから漢字に切り替わる位置等）で語の単位に分解したりする方法がある。類似度計算部3は、同義語辞書を検索して、分解した語が同義語辞書に同義語として登録されていると、その同義語に対応づけられている代表語に変換する。

【0039】

40

類似度計算部3は、省略可能語辞書を検索して、代表語に変換された語が省略可能語として登録されていると、同じグループIDが付されている省略可能語を、データベース2に登録されているレコードの語に追加する。そして、類似度計算部3は、代表語に変換され、省略可能語が追加された各レコード間の類似度を計算する。

【0040】

類似度計算部3が各レコード間の類似度を計算する方法は、例えば、各情報の対応するフィールドの語の文字を先頭から1文字ずつ比較していき、合致すればその文字の類似度を1とし、合致しなければその文字の類似度を0とする。そして、例えば、各文字の類似度を合計した数を、語の文字数で割った商（すなわち、0から1の間で正規化した値）をそのフィールドの類似度とする。なお、各情報のフィールドの語の文字数が異なっている

50

場合は、各フィールドを構成する語のうち最も多い文字数で、各文字の類似度を合計した数を割った商をそのフィールドの類似度とする。

【0041】

そして、類似度計算部3は、各フィールドの類似度に、各フィールドごとの所定の重みの値を乗じた積を合計した数を、フィールドの数で割った商（すなわち、0から1の間で正規化した値）を、レコードの類似度として計算する。

【0042】

本発明の類似度の計算方法は、上述した方法に限定されるものではなく、編集距離を用いる方法等の、他の方法を用いてもよい。なお、類似度計算部3は、データベース情報入力部4に入力された情報にもとづいて、各フィールドおよび各レコードの類似度を計算し

10

【0043】

重複候補抽出部6は、類似度計算部3が計算した類似度が、所定の閾値以上であるレコードの組を、重複レコード候補として抽出する。なお、重複候補抽出部6は、他の方法を用いて、重複レコード候補を抽出してもよい。

【0044】

重複レコード検出システム20は、コンピュータに、複数の情報からなる複数のレコードを保持するデータベースに登録されている情報に用いられている語を、語に対応する変換語を記憶する変換語辞書5が記憶している変換語に変換して、レコード間の類似度を計算させる類似度計算処理と、類似度計算処理で計算した類似度が、所定の値以上であるレ

20

【0045】

次に、本発明の第1の実施の形態の動作を、具体例を挙げて図面を参照して説明する。図2の例に示すデータベース2に登録されている情報から、重複レコード候補を抽出する。図4は、本発明の第1の実施の形態の動作を説明するフローチャートである。

【0046】

まず、類似度計算部3が、データベース2に登録されている情報を読み込む（ステップS101）。レコードIDが「001」のレコードの名称のフィールドは、「エヌイーシー奈良支店」である。類似度計算部3は、「エヌイーシー奈良支店」に形態素解析等を行って語に分解する（ステップS102）。具体的には、「エヌイーシー」と「奈良」と「支店」とに分解する。

30

【0047】

なお、ここでは、データベース情報入力部4に入力された情報が、名称フィールド、および住所フィールドの類似度計算に、変換語辞書5を用いることを示していたものとする。すると、類似度計算部3は、変換語辞書5を参照して、分解した語が同義語であれば代表語に変換する（ステップS103）。図3を参照すると、「エヌイーシー」が同義語であるので、「エヌイーシー」を代表語である「NEC」に変換する。「奈良」および「支店」は同義語辞書に登録されていないので変換を行わない。すると、レコードID「001」の語は、「NEC」、「奈良」、および「支店」である。

40

【0048】

次に、類似度計算部3は、変換語辞書5を参照して、省略可能語があれば、同じグループIDの省略可能語を追加する（ステップS104）。図3を参照すると、「NEC」と「日電」とが同じグループIDの省略可能語であるので、「日電」を追加する。すると、レコードID「001」の語は、「NEC」、「日電」、「奈良」、および「支店」である。

【0049】

類似度計算部3は、分解した語を結合する（ステップS105）。すると、レコードID「001」のレコードの名称のフィールドは、「NEC日電奈良支店」および「日電NEC奈良支店」に変換される。

50

【 0 0 5 0 】

類似度計算部 3 は、上述したステップ S 1 0 1 からステップ S 1 0 5 の動作を、レコード ID 「 0 0 2 」およびレコード ID 「 0 0 3 」に対しても行う。

【 0 0 5 1 】

具体的には、レコード ID が 「 0 0 2 」のレコードの名称のフィールドは、「日電奈良支店」である。類似度計算部 3 は、「日電奈良支店」に形態素解析を行って語に分解する。具体的には、「日電」と「奈良」と「支店」とに分解する。

【 0 0 5 2 】

類似度計算部 3 は、変換語辞書 5 を参照して、分解した各語が同義語であれば代表語に変換する。図 3 を参照すると、「日電」は代表語であり、「奈良」および「支店」は同義語辞書に登録されていないので変換を行わない。

10

【 0 0 5 3 】

次に、類似度計算部 3 は、変換語辞書 5 を参照して、省略可能語があれば、同じグループ ID の省略可能語を追加する。図 3 を参照すると、「NEC」と「日電」とが同じグループ ID の省略可能語であるので、「NEC」を追加する。すると、レコード ID 「 0 0 2 」の語は、「NEC」、「日電」、「奈良」、および「支店」である。

【 0 0 5 4 】

類似度計算部 3 は、分解した語を結合する。すると、レコード ID 「 0 0 2 」のレコードの名称のフィールドは、「NEC日電奈良支店」および「日電NEC奈良支店」に変換される。

20

【 0 0 5 5 】

同様に、レコード ID が 「 0 0 3 」のレコードの名称のフィールドは、「NEC奈良支店」である。類似度計算部 3 は、「NEC奈良支店」に形態素解析を行って語に分解する。具体的には、「NEC」と「奈良」と「支店」とに分解する。

【 0 0 5 6 】

類似度計算部 3 は、変換語辞書 5 を参照して、分解した各語が同義語であれば代表語に変換する。図 3 を参照すると、「NEC」が代表語であり、「奈良」および「支店」は同義語辞書に登録されていないので変換を行わない。

【 0 0 5 7 】

次に、類似度計算部 3 は、変換語辞書 5 を参照して、省略可能語があれば、同じグループ ID の省略可能語を追加する。図 3 を参照すると、「NEC」と「日電」とが同じグループ ID の省略可能語であるので、「日電」を追加する。すると、レコード ID 「 0 0 6 」の語は、「NEC」、「日電」、「奈良」、および「支店」である。

30

【 0 0 5 8 】

類似度計算部 3 は、分解した語を結合する。すると、レコード ID 「 0 0 3 」のレコードの名称のフィールドは、「NEC日電奈良支店」および「日電NEC奈良支店」に変換される。

【 0 0 5 9 】

次に、類似度計算部 3 は、変換したレコード ID 「 0 0 1 」、 「 0 0 2 」および「 0 0 3 」の名称のフィールドの相互の類似度を計算する(ステップ S 1 0 6)。

40

【 0 0 6 0 】

まず、レコード ID 「 0 0 1 」の名称のフィールドと、レコード ID 「 0 0 2 」の名称のフィールドとの類似度を計算する。レコード ID 「 0 0 1 」の変換後の名称のフィールドは、「NEC日電奈良支店」と、「日電NEC奈良支店」とであり、レコード ID 「 0 0 2 」の変換後の名称のフィールドは、「NEC日電奈良支店」と、「日電NEC奈良支店」とである。レコード ID 「 0 0 1 」の「NEC日電奈良支店」と、レコード ID 「 0 0 2 」の「NEC日電奈良支店」とは、9文字中9文字が合致するので、 $9 \times 1 \div 9 = 1$ となり、類似度は1である。

【 0 0 6 1 】

同様に、レコード ID 「 0 0 1 」の「日電NEC奈良支店」と、レコード ID 「 0 0 2 」

50

」の「日電NEC奈良支店」とは、9文字中9文字が合致するので、 $9 \times 1 \div 9 = 1$ となり、類似度は1である。

【0062】

また、レコードID「001」の「NEC日電奈良支店」と、レコードID「002」の「日電NEC奈良支店」とは、9文字中4文字が合致するので、 $4 \times 1 \div 9 = 0.44$ （小数点3桁目四捨五入）となり、類似度は0.44である。

【0063】

同様に、また、レコードID「001」の「日電NEC奈良支店」と、レコードID「002」の「NEC日電奈良支店」とは、9文字中4文字が合致するので、 $4 \times 1 \div 9 = 0.44$ （小数点3桁目四捨五入）となり、類似度は0.44である。

10

【0064】

ここで、類似度計算部3は、最も類似度が高い値を採用することとする。すると、レコードID「001」の名称のフィールドと、レコードID「002」の名称のフィールドとの類似度は1である。

【0065】

同様に、レコードID「001」の名称のフィールドと、レコードID「003」の名称のフィールドとの類似度を計算すると、類似度は1となる。また、レコードID「002」の名称のフィールドと、レコードID「003」の名称のフィールドとの類似度を計算すると、類似度は1となる。

【0066】

20

次に、類似度計算部3は、レコードID「001」、「002」および「003」の住所のフィールドの相互の類似度を計算する（ステップS107）。

【0067】

レコードID「001」の住所のフィールドは「1の1」であり、レコードID「002」の住所のフィールドは「1-1」であり、レコードID「003」の住所のフィールドは「1-1」である。

【0068】

レコードID「001」の住所のフィールド「1の1」と、レコードID「002」の住所のフィールド「1-1」とは、6文字中5文字が合致するので、 $5 \times 1 \div 6 = 0.83$ （小数点3桁目四捨五入）となり、類似度は0.83である。

30

【0069】

レコードID「001」の住所のフィールド「1の1」と、レコードID「003」の住所のフィールド「1-1」とは、6文字中5文字が合致するので、 $5 \times 1 \div 6 = 0.83$ （小数点3桁目四捨五入）となり、類似度は0.83である。

【0070】

レコードID「002」の住所のフィールド「1-1」と、レコードID「003」の住所のフィールド「1-1」とは、6文字中6文字が合致するので、 $6 \times 1 \div 6 = 1$ となり、類似度は1である。

【0071】

次に、類似度計算部3は、レコードID「001」、「002」および「003」の電話番号のフィールドの相互の類似度を計算する（ステップS108）。

40

【0072】

レコードID「001」の電話番号のフィールドは「000-111-1234」であり、レコードID「002」の電話番号のフィールドは「000-111-1235」であり、レコードID「003」の電話番号のフィールドは「000-111-1234」である。

【0073】

レコードID「001」の電話番号のフィールド「000-111-1234」と、レコードID「002」の電話番号のフィールド「000-111-1235」とは、12文字中11文字が合致するので、 $11 \times 1 \div 12 = 0.92$ （小数点3桁目四捨五入）と

50

なり、類似度は0.92である。

【0074】

レコードID「001」の電話番号のフィールド「000-111-1234」と、レコードID「003」の電話番号のフィールド「000-111-1234」とは、12文字中12文字が合致するので、 $12 \times 1 \div 12 = 1$ となり、類似度は1である。

【0075】

レコードID「002」の電話番号のフィールド「000-111-1235」と、レコードID「003」の電話番号のフィールド「000-111-1234」とは、12文字中11文字が合致するので、 $11 \times 1 \div 12 = 0.92$ （小数点3桁目四捨五入）となり、類似度は0.92である。

10

【0076】

類似度計算部3は、レコードID「001」、「002」および「003」の各フィールドの相互の類似度を、各フィールドの類似度に重みの値を乗じた積を合計した数を、フィールドの数で割った商を、各レコード間の類似度として計算する（ステップS109）。なお、ここでは、データベース情報入力部4に入力された情報が、名称フィールド、住所フィールド、および電話番号フィールドの重みの比が、1:1:1であることを示していたものとする。

【0077】

すると、レコードID「001」とレコードID「002」との類似度は、 $(1 \times 1 + 0.83 \times 1 + 0.92 \times 1) \div 3 = 0.92$ （小数点3桁目四捨五入）となる。

20

【0078】

また、レコードID「001」とレコードID「003」との類似度は、 $(1 \times 1 + 0.83 \times 1 + 1 \times 1) \div 3 = 0.94$ （小数点3桁目四捨五入）となる。

【0079】

レコードID「002」とレコードID「003」との類似度は、 $(1 \times 1 + 1 \times 1 + 0.92 \times 1) \div 3 = 0.97$ （小数点3桁目四捨五入）となる。

【0080】

類似度計算部3は、計算した各レコードの組の類似度と、類似度を計算したレコードの組とを重複候補抽出部6に出力する。重複候補抽出部6は、類似度計算部3が計算した類似度が、所定の閾値以上である各レコードを、重複レコード候補として抽出する（ステップS110）。ここで、所定の閾値を0.90とすると、重複候補抽出部6は、レコードID「001」、レコードID「002」、およびレコードID「003」を重複レコード候補として抽出する。

30

【0081】

表示部（図示せず）は、重複候補抽出部6が抽出した各レコードを表示する（ステップS111）。

【0082】

以上に述べたように、この実施の形態によれば、同義語や省略可能語による表記の差異があっても、重複する情報をデータベース2から抽出することができる。

【0083】

実施の形態2.

40

本発明の第2の実施の形態について、図面を参照して説明する。図5は、本発明の第2の実施の形態の一構成例を示すブロック図である。

【0084】

本発明の第2の実施の形態の構成は、第1の実施の形態の構成に、重複候補抽出部6が抽出した重複レコード候補が、重複レコードであるか否かを判定するルールである重複判定ルールを記憶する重複判定ルール記憶部8、重複判定ルール記憶部8が記憶しているルールにもとづいて、重複候補抽出部6が抽出した重複レコード候補が重複レコードであるか否かを判定する重複レコード判定部7、および重複レコード判定部7が重複レコードであると判定したレコードを記憶する重複レコードデータベース9を加えたものであり、そ

50

他の構成要素は第 1 の実施の形態と同様なため、その他の構成要素には図 1 と同じ符号を付し、説明を省略する。

【 0 0 8 5 】

図 6 は、重複判定ルール例を示す説明図である。図 6 に示した例によると、重複判定ルールは、例えば、レコード相互の類似度が特定の値を超えているならば、それらを重複レコードとみなす、というルールや、レコード相互の類似度が特定の値以下であれば、それらを重複レコードではないとみなす、というルールや、いずれかのフィールドの類似度が特定の値以下であれば、それらを重複レコードではないとみなす、というルールや、あるフィールドの類似度が所定の値以上であり、かつ、他のあるフィールドの類似度が所定の値以上であれば、それらを重複レコードとみなす、等である。

10

【 0 0 8 6 】

図 7 は、重複判定ルール記憶部 8 が記憶している重複判定ルール例を示す説明図である。図 7 の例に示すように、重複判定ルールは、それぞれ条件部分（図 7 における I F 以下の部分）と、結論部分（図 7 における T H E N 以下の部分）とで構成される。

【 0 0 8 7 】

そして、条件部分には、レコードの組の類似度の値や、フィールドの組の類似度の値が、ある値よりも大きい、小さい、以上、または以下等の条件を、A N D、O R、および N O T で組み合わせて記述する。

【 0 0 8 8 】

また、結論部分には、条件部分の記述されている条件に合致するレコードの組を、重複レコードであると記述したり、重複レコードではないと記述したりする。また、結論部分において、条件部分のネスト（入れ子）を記述してもよい。

20

【 0 0 8 9 】

図 7 の（ 1 ）式に示す例では、レコードの組の類似度の値が 1 であれば、重複レコードであるというルールを記述している。また、図 7 の（ 2 ）式に示す例では、住所フィールドの組の類似度が 0 . 9 を超えていて、かつ、電話番号フィールドの組の類似度が 0 . 9 を超えていた場合、名称フィールドの組の類似度が 0 . 9 を超えていれば、レコードの組は重複レコードであり、名称フィールドの組の類似度が 0 . 9 以下であれば、レコードの組を重複レコードではないというルールを記述している。

【 0 0 9 0 】

類似度計算部 3 は、計算した各フィールドの組の類似度と、各レコードの組の類似度とを重複候補抽出部 6 に出力する。重複候補抽出部 6 は、類似度計算部 3 が計算した類似度が、所定の閾値以上である各レコードを、重複レコード候補として抽出し、重複レコード候補の各フィールドの組の類似度と、各レコードの組の類似度とを重複レコード判定部 7 に出力する。

30

【 0 0 9 1 】

本発明の第 2 の実施の形態の動作を、具体例を挙げて説明する。まず、図 2 の例に示したレコード I D 「 0 0 1 」、レコード I D 「 0 0 2 」およびレコード I D 「 0 0 3 」が重複レコードであるか否かを判定する場合を例に説明する。

【 0 0 9 2 】

類似度計算部 3 が、各フィールドの組の類似度と、各レコードの組の類似度とを計算するまでの動作は、第 1 の実施の形態における動作と同様なため、説明を省略する。

40

【 0 0 9 3 】

類似度計算部 3 は、計算した各フィールドの組の類似度と、各レコードの組の類似度とを重複候補抽出部 6 に出力する。重複候補抽出部 6 は、類似度計算部 3 が計算した類似度が、所定の閾値以上である各レコードを重複レコード候補として抽出し、重複レコード候補の各フィールドの組の類似度と、各レコードの組の類似度とを重複レコード判定部 7 に出力する。ここで、所定の閾値を 0 . 9 とすると、重複候補抽出部 6 は、レコード I D 「 0 0 1 」、レコード I D 「 0 0 2 」、およびレコード I D 「 0 0 3 」を重複レコード候補として抽出する。

50

【 0 0 9 4 】

重複レコード判定部 7 は、レコード ID 「 0 0 1 」 とレコード ID 「 0 0 2 」 との類似度が 0 . 9 2 であるので、図 7 の例に示す式 (1) の条件部分 (レコードの組の類似度の値が 1) に合致しないので、レコード ID 「 0 0 1 」 とレコード ID 「 0 0 2 」 との重複レコードの判定に、式 (1) を適用しない。

【 0 0 9 5 】

重複レコード判定部 7 は、レコード ID 「 0 0 1 」 の住所フィールドと、レコード ID 「 0 0 2 」 の住所フィールドとの類似度が 1 であるが、レコード ID 「 0 0 1 」 の電話番号フィールドと、レコード ID 「 0 0 2 」 の電話番号フィールドとの類似度が 0 . 8 3 であるので、図 7 の例に示す式 (2) の条件部分 (住所フィールドの組の類似度が 0 . 9 を超えていて、かつ、電話番号フィールドの組の類似度が 0 . 9 を超えていた場合) に合致しないので、レコード ID 「 0 0 1 」 とレコード ID 「 0 0 2 」 との重複レコードの判定に、式 (2) を適用しない。

10

【 0 0 9 6 】

また、重複レコード判定部 7 は、レコード ID 「 0 0 1 」 とレコード ID 「 0 0 3 」 との類似度が 0 . 9 4 であるので、図 7 の例に示す式 (1) の条件部分 (レコードの組の類似度の値が 1) に合致しないので、レコード ID 「 0 0 1 」 とレコード ID 「 0 0 3 」 との重複レコードの判定に、式 (1) を適用しない。

【 0 0 9 7 】

重複レコード判定部 7 は、レコード ID 「 0 0 1 」 の住所フィールドと、レコード ID 「 0 0 3 」 の住所フィールドとの類似度が 1 であるが、レコード ID 「 0 0 1 」 の電話番号フィールドと、レコード ID 「 0 0 3 」 の電話番号フィールドとの類似度が 0 . 8 3 であるので、図 7 の例に示す式 (2) の条件部分 (住所フィールドの組の類似度が 0 . 9 を超えていて、かつ、電話番号フィールドの組の類似度が 0 . 9 を超えていた場合) に合致しないので、レコード ID 「 0 0 1 」 とレコード ID 「 0 0 3 」 との重複レコードの判定に、式 (2) を適用しない。

20

【 0 0 9 8 】

重複レコード判定部 7 は、レコード ID 「 0 0 2 」 とレコード ID 「 0 0 3 」 との類似度が 0 . 9 7 であるので、図 7 の例に示す式 (1) の条件部分 (レコードの組の類似度の値が 1) に合致しないので、レコード ID 「 0 0 2 」 とレコード ID 「 0 0 3 」 との重複レコードの判定に、式 (1) を適用しない。

30

【 0 0 9 9 】

重複レコード判定部 7 は、レコード ID 「 0 0 2 」 の住所フィールドと、レコード ID 「 0 0 3 」 の住所フィールドとの類似度が 1 であって、レコード ID 「 0 0 2 」 の電話番号フィールドと、レコード ID 「 0 0 3 」 の電話番号フィールドとの類似度が 1 であるので、図 7 の例に示す式 (2) の条件部分 (住所フィールドの組の類似度が 0 . 9 を超えていて、かつ、電話番号フィールドの組の類似度が 0 . 9 を超えていた場合) に合致する。また、レコード ID 「 0 0 2 」 の名称フィールドと、レコード ID 「 0 0 3 」 の名称フィールドとの類似度が 0 . 9 2 であるので、式 (2) の結果部分における条件部分 (名称フィールドの組の類似度が 0 . 9 を超えている) に合致するので、レコード ID 「 0 0 2 」 とレコード ID 「 0 0 3 」 とが重複レコードであると判定する。

40

【 0 1 0 0 】

重複レコード判定部 7 は、重複レコードであると判定した各レコードを、重複レコードデータベース 9 に記憶させる。

【 0 1 0 1 】

なお、重複レコード判定部 7 は、重複判定ルール記憶部 8 が記憶している重複判定ルールを適用しなかったレコードの組を、表示部に表示させてもよい。すると、使用者が重複レコードであるか否かを判定することができる。

【 0 1 0 2 】

以上に述べたように、この実施の形態によれば、重複レコード判定部 7 が、予め重複判

50

定ルール記憶部 8 が記憶している重複判定ルールにもとづいて、各レコードの組が重複レコードであるか否かを自動的に判定することができる。

【0103】

また、重複レコードデータベース 9 が、重複レコード判定部 7 が重複レコードであると判定したレコードの組を記憶するため、使用者は、重複レコード判定部 7 が重複レコードであると判定したレコードの組を確認することができる。

【0104】

実施の形態 3 .

本発明の第 3 の実施の形態を、図面を参照して説明する。図 8 は、本発明の第 3 の実施の形態の一構成例を示すブロック図である。

10

【0105】

本発明の第 3 の実施の形態の構成は、第 2 の実施の形態の構成に、重複レコード判定部 7 が重複レコードであると判定したレコードの組から変換語の候補を抽出して変換語辞書 5 に登録する変換語候補抽出部 10 を加えた点が第 2 の実施の形態の構成と異なり、その他の点は第 2 の実施の形態の構成と同様である。そのため、第 2 の実施の形態と同様な構成要素には、図 5 と同じ符号を付し、説明を省略する。

【0106】

変換語候補抽出部 10 は、重複レコード判定部 7 が重複レコードであると判定したレコードの組を比較して、異なる部分に、例えば、形態素解析等を行って、重複レコードの組における異なる部分の語の組を同義語候補の組として抽出する。

20

【0107】

なお、変換語候補抽出部 10 は、重複レコード判定部 7 が重複レコードであると判定したレコードの組を比較して、一方のレコードが、他方のレコードに含まれる場合には、一方のレコードと他方のレコードとの異なる部分と、共通する部分との組を省略可能語候補の組として抽出する。

【0108】

また、変換語候補抽出部 10 は、抽出した同義語候補の組が、他の一のレコードに含まれる場合は、抽出した同義語候補の組を省略可能語候補の組とする。

【0109】

変換語候補抽出部 10 は、抽出した同義語候補の組のうち、他の同義語候補や省略可能語候補の組み合わせで構成される同義語候補の組を、同義語候補の組から除外する。

30

【0110】

また、変換語候補抽出部 10 は、省略可能語候補の組のうち、他の同義語候補や省略可能語候補に含まれる省略可能語候補の組を、省略可能語候補の組から除外する。

【0111】

変換語候補抽出部 10 は、変換語辞書 5 を参照して、変換語候補、および省略可能語候補の組のうち、変換語辞書 5 に登録されている語以外の語を変換語辞書 5 に登録する。

【0112】

次に、この実施の形態において、重複レコードから同義語候補および省略可能語候補を名称フィールドから抽出する際の動作を、具体例を挙げて図面を参照して説明する。図 9 は、本発明の第 3 の実施の形態の動作を説明するフローチャートである。図 10 は、重複レコード判定部 7 が重複レコードであると判定したレコードの組の例を示す説明図である。

40

【0113】

変換語候補抽出部 10 は、重複レコード判定部 7 が重複レコードであると判定したレコードの組を比較して、異なる部分に、例えば、形態素解析等を行って、重複レコードの組における異なる部分の語を同義語候補として抽出し、一方のレコードが、他方のレコードに含まれる場合には、一方のレコードと他方のレコードとの異なる部分と、共通する部分とを省略可能語候補として抽出する（ステップ S 3 0 1）。

【0114】

50

具体的には、変換語候補抽出部10は、図10の例に示したレコードID「001」とレコードID「002」とを比較して、レコードID「001」と、レコードID「002」とで異なる部分である「日電NEC」と「エヌイーシー」とを同義語候補の組として抽出する。

【0115】

変換語候補抽出部10は、図10の例に示したレコードID「001」とレコードID「003」とを比較して、レコードID「001」と、レコードID「003」とで異なる部分である「NEC」と「エヌイーシー」とを同義語候補の組として抽出する。

【0116】

変換語候補抽出部10は、図10の例に示したレコードID「001」とレコードID「004」とを比較すると、レコードID「004」が、レコードID「001」に含まれるので、異なる部分である「日電」と、共通する部分である「NEC奈良支店」とを省略可能語候補の組として抽出する。

10

【0117】

変換語候補抽出部10は、図10の例に示したレコードID「001」とレコードID「005」とを比較すると、レコードID「005」が、レコードID「001」に含まれるので、異なる部分である「NEC」と、共通する部分である「日電奈良支店」とを省略可能語候補の組として抽出する。

【0118】

変換語候補抽出部10は、図10の例に示したレコードID「002」とレコードID「003」とを比較すると、レコードID「002」が、レコードID「003」に含まれるので、異なる部分である「日電」と、共通する部分である「エヌイーシー奈良支店」とを省略可能語候補の組として抽出する。

20

【0119】

変換語候補抽出部10は、図10の例に示したレコードID「002」とレコードID「004」とを比較して、レコードID「002」と、レコードID「004」とで異なる部分である「エヌイーシー」と「NEC」とを同義語候補の組として抽出する。

【0120】

変換語候補抽出部10は、図10の例に示したレコードID「002」とレコードID「005」とを比較して、レコードID「002」と、レコードID「005」とで異なる部分である「エヌイーシー」と「日電」とを同義語候補の組として抽出する。

30

【0121】

変換語候補抽出部10は、図10の例に示したレコードID「003」とレコードID「004」とを比較して、レコードID「003」と、レコードID「004」とで異なる部分である「日電エヌイーシー」と「NEC」とを同義語候補の組として抽出する。

【0122】

変換語候補抽出部10は、図10の例に示したレコードID「003」とレコードID「005」とを比較すると、レコードID「005」が、レコードID「003」に含まれるので、異なる部分である「エヌイーシー」と、共通する部分である「日電奈良支店」とを省略可能語候補の組として抽出する。

40

【0123】

変換語候補抽出部10は、図10の例に示したレコードID「004」とレコードID「005」とを比較して、レコードID「004」と、レコードID「005」とで異なる部分である「NEC」と「日電」とを同義語候補の組として抽出する。

【0124】

次に、変換語候補抽出部10は、抽出した同義語候補の組が、他の一のレコードに含まれる場合は、抽出した同義語候補の組を省略可能語候補の組とする(ステップS302)。

【0125】

具体的には、変換語候補抽出部10は、レコードID「002」とレコードID「00

50

5」とを比較して抽出した同義語候補の組である「エヌイーシー」と「日電」とが、レコードID「003」の「日電エヌイーシー奈良支店」に含まれるので、同義語候補の組である「エヌイーシー」と「日電」とを省略可能語候補の組とする。

【0126】

また、変換語候補抽出部10は、レコードID「004」とレコードID「005」とを比較して抽出した同義語候補の組である「NEC」と「日電」とが、レコードID「001」の「日電NEC奈良支店」に含まれるので、同義語候補の組である「NEC」と「日電」と省略可能語候補の組とする。

【0127】

次に、変換語候補抽出部10は、抽出した同義語候補の組のうち、他の同義語候補や省略可能語候補の組み合わせで構成される同義語候補の組を、同義語候補の組から除外する(ステップS303)。

10

【0128】

具体的には、変換語候補抽出部10は、レコードID「004」とレコードID「005」とを比較して抽出した同義語候補の組である「NEC」と「日電」とを組み合わせると、レコードID「001」とレコードID「002」とを比較して同義語候補として抽出した「日電NEC」を構成するので、レコードID「004」とレコードID「005」とを比較して抽出した同義語候補の組である「NEC」と「日電」とを、同義語候補の組から除外する。

20

【0129】

変換語候補抽出部10は、レコードID「002」とレコードID「005」とを比較して抽出した同義語候補の組である「エヌイーシー」と「日電」とを組み合わせると、レコードID「003」とレコードID「004」とを比較して同義語候補として抽出した「日電エヌイーシー」を構成するので、レコードID「002」とレコードID「005」とを比較して抽出した同義語候補の組である「エヌイーシー」と「日電」とを、同義語候補の組から除外する。

30

【0130】

変換語候補抽出部10は、省略可能語候補の組のうち、他の同義語候補や省略可能語候補に含まれる省略可能語候補の組を、省略可能語候補の組から除外する(ステップS304)。

30

【0131】

具体的には、変換語候補抽出部10は、レコードID「001」とレコードID「004」とを比較して抽出した省略可能語候補の組である「日電」と「NEC奈良支店」とは、レコードID「004」とレコードID「005」とを比較して抽出した同義語候補の組である「NEC」と「日電」とを含むので省略可能語候補から除外する。

【0132】

変換語候補抽出部10は、レコードID「001」とレコードID「005」とを比較して抽出した省略可能語候補の組である「NEC」と「日電奈良支店」とは、レコードID「004」とレコードID「005」とを比較して抽出した同義語候補の組である「NEC」と「日電」とを含むので省略可能語候補から除外する。

40

【0133】

変換語候補抽出部10は、レコードID「002」とレコードID「003」とを比較して抽出した省略可能語候補の組である「日電」と「エヌイーシー奈良支店」とは、レコードID「002」とレコードID「005」とを比較して抽出した同義語候補の組である「エヌイーシー」と「日電」とを含むので省略可能語候補から除外する。

【0134】

変換語候補抽出部10は、レコードID「003」とレコードID「005」とを比較して抽出した省略可能語候補の組である「エヌイーシー」と「日電奈良支店」とは、レコードID「002」とレコードID「005」とを比較して抽出した同義語候補の組である「エヌイーシー」と「日電」とを含むので省略可能語候補から除外する。

50

【 0 1 3 5 】

すると、レコードID「001」とレコードID「003」とを比較して抽出された同義語候補の組、およびレコードID「002」とレコードID「004」とを比較して抽出された同義語候補の組である「NEC」と「エヌイーシー」とが同義語候補の組となる。

【 0 1 3 6 】

また、レコードID「002」とレコードID「005」とを比較して抽出された同義語候補の組、およびレコードID「004」とレコードID「005」とを比較して抽出された同義語候補の組である「エヌイーシー」と「日電」とが、省略可能語候補に変更されて省略可能語候補の組となる。

10

【 0 1 3 7 】

変換語候補抽出部10は、変換語辞書5を参照して、変換語候補、および省略可能語候補の組のうち、変換語辞書5に登録されている語以外の語を変換語辞書5に登録する(ステップS305)。

【 0 1 3 8 】

なお、変換語候補抽出部10は、同義語候補の組である「NEC」と「エヌイーシー」とのいずれかを代表語として、変換語辞書5に登録する。変換語候補抽出部10は、例えば、50音順のや、アルファベット順の早い方の語や、文字数の少ない方の語を代表語として、同義語候補の組の語を変換語辞書5に登録する。

20

【 0 1 3 9 】

このとき、同義語候補の組の語のいずれかが既に代表語として変換語辞書5に登録されていた場合、変換語候補抽出部10は、同義語候補の組の他の語を、その代表語の同義語として変換語辞書5に登録する。

【 0 1 4 0 】

なお、同義語候補の組の語のすべてが既に代表語として変換語辞書5に登録されていた場合、変換語候補抽出部10は、いずれか1つの語を代表語として変換語辞書5に登録し、同義語候補の組の他の語を、その代表語の同義語として変換語辞書5に登録する。このとき、変換語候補抽出部10は、既に代表語として変換語辞書5に登録されていた語の同義語を、代表語として変換語辞書5に登録した語の同義語として、変換語辞書5に登録する。

30

【 0 1 4 1 】

変換語候補抽出部10は、省略可能語候補の組が変換語辞書5に登録されていなかった場合、新たにグループIDを決定して省略可能語候補として変換語辞書5に登録する。なお、新たなグループIDは、例えば、既に登録済みのグループIDの最大値に1を加えた値とする。

【 0 1 4 2 】

変換語候補抽出部10は、省略可能語候補の組のうち、いずれかが既に省略可能語候補として変換語辞書5に登録されていた場合、他の省略可能語候補を、既に変換語辞書5に登録されている省略可能語候補と同じグループIDで変換語辞書5に登録する。

【 0 1 4 3 】

変換語候補抽出部10は、省略可能語候補の組がすべて既に省略可能語候補として同じグループIDで変換語辞書5に登録されていた場合は、変換語辞書5に登録する動作を行わない。変換語候補抽出部10は、省略可能語候補の組を構成する省略可能語候補のそれぞれが、異なるグループIDで既に省略可能語候補として変換語辞書5に登録されていた場合、グループIDの値を比較して、グループIDの値が大きい方の省略可能語候補のグループIDの値を、グループIDが小さい方の値に変更する。

40

【 0 1 4 4 】

なお、変換語候補抽出部10は、同義語候補または省略可能語候補を表示部に表示して、同義語候補または省略可能語候補として変換語辞書5に登録するか否かを使用者に選択させてもよい。

50

【 0 1 4 5 】

以上に述べたように、この実施の形態によれば、重複レコード判定部 7 が重複レコードとして判定したレコードに含まれる語のうち、所定の条件に合致する語を、同義語候補または省略可能語候補として変換語辞書 5 に登録することができる。

【 0 1 4 6 】

実施の形態 4 .

本発明の第 4 の実施の形態を、図面を参照して説明する。図 1 1 は、本発明の第 4 の実施の形態の一構成例を示すブロック図である。

【 0 1 4 7 】

本発明の第 4 の実施の形態の構成は、第 3 の実施の形態の構成に、同義語候補、および省略可能語候補を変換語辞書 5 に登録するか否か、および重複レコード候補を重複レコードデータベース 9 に登録するか否かを使用者に確認する入出力部 1 1 を加えたものであり、その他の構成要素は第 3 の実施の形態と同様のため、その他の構成要素には図 8 と同じ符号を付し、説明を省略する。

10

【 0 1 4 8 】

入出力部 1 1 は、表示部である液晶ディスプレイ等と、入力手段であるキーボード等によって実現される。

【 0 1 4 9 】

次に、本発明の第 4 の実施の形態の動作を、図面を参照して説明する。図 1 2 は、本発明の第 4 の実施の形態の動作を説明するフローチャートである。

20

【 0 1 5 0 】

まず、類似度計算部 3 が、データベース 2 に登録されている情報を読み込む（ステップ S 4 0 1 ）。

【 0 1 5 1 】

類似度計算部 3 は、第 1 の実施の形態における動作と同様の動作を行い、各レコードの組の類似度を計算し、計算した各レコードの組の類似度と、類似度を計算したレコードの組とを重複候補抽出部 6 に出力する（ステップ S 4 0 2 ）。

【 0 1 5 2 】

重複候補抽出部 6 は、例えば、第 1 の実施の形態における動作と同様の動作を行い、重複レコード候補をデータベース 2 から抽出する（ステップ S 4 0 3 ）。

30

【 0 1 5 3 】

重複レコード判定部 7 は、第 2 の実施の形態における動作と同様の動作を行い、重複候補抽出部 6 が抽出した重複レコード候補が重複レコードであるか否かを判定し（ステップ S 4 0 4 ）、重複レコードであると判定した重複レコード候補を重複レコードデータベース 9 に記憶させる。

【 0 1 5 4 】

変換語候補抽出部 1 0 は、第 3 の実施の形態における動作と同様の動作を行い、重複レコード判定部 7 が重複レコードであると判定したレコードから、同義語候補および省略可能語候補（以下、単に変換語候補という）を抽出する（ステップ S 4 0 5 ）。

【 0 1 5 5 】

変換語候補抽出部 1 0 は、変換語候補を抽出すると、抽出した変換語候補を入出力部 1 1 に出力する。入出力部 1 1 は、変換語候補抽出部 1 0 が入力した変換語候補を表示し、使用者に変換語候補が変換語であるか否かを判定させる。

40

【 0 1 5 6 】

使用者が、入出力部 1 1 を操作して変換語候補が変換語であると判定すると、変換語候補抽出部 1 0 は、変換語候補を変換語辞書 5 に登録する（ステップ S 4 0 6 ）。

【 0 1 5 7 】

変換語候補抽出部 1 0 が、変換語候補を変換語辞書 5 に登録すると、重複候補抽出部 6 は、第 1 の実施の形態における動作と同様の動作を行い、重複レコード候補をデータベース 2 から抽出する（ステップ S 4 0 3 ）。新たな変換語が変換語辞書 5 に登録されると、

50

新たに重複レコードと判定されるレコードが発生する可能性があるからである。

【0158】

変換語候補抽出部10が変換語候補の抽出を終了したり、使用者が、入出力部11を操作して変換語候補が変換語であるか否かの判定を拒否したりすると、重複レコード判定部7は、重複レコードであると判定しなかった重複レコード候補を入出力部11に出力する。入出力部11は、重複レコード判定部7が入力した重複レコード候補を表示し、使用者に重複レコード候補が重複レコードであるか否かを判定させる(ステップS407)。なお、重複レコード判定部7は、重複レコードであると判定しなかった重複レコード候補のうち、重複する可能性のあるレコードの組み合わせの数が多い順番で、重複レコード候補を入出力部11に出力してもよい。

10

【0159】

使用者が、入出力部11を操作して重複レコード候補が重複レコードであると判定すると、重複レコード判定部7は、重複レコード候補を重複レコードであると判定し、重複レコードデータベース9に登録する(ステップS408)。

【0160】

重複レコード判定部7が、重複レコード候補を重複レコードデータベース9に登録すると、変換語候補抽出部10は、第3の実施の形態における動作と同様の動作を行い、重複レコード判定部7が重複レコードであると判定したレコードから、変換語候補を抽出する(ステップS405)。重複レコードが増加すると、増加した重複レコードから新たに交換語候補が抽出される可能性があるからである。

20

【0161】

重複レコード判定部7が、重複レコードであると判定しなかった重複レコード候補をすべて入出力部11に出力したり、使用者が、入出力部11を操作して重複レコード候補の判定を拒否したりすると、変換語候補抽出部10は、第3の実施の形態における動作と同様の動作を行い、重複レコード判定部7が重複レコードであると判定したレコードから変換語候補を抽出する(ステップS405)。

【0162】

変換語候補抽出部10が、第3の実施の形態における動作と同様の動作を行い、重複レコード判定部7が重複レコードであると判定したレコードから変換語候補を抽出する動作を終了すると、すべての構成要素は動作を終了する。

30

【0163】

以上に述べたように、この実施の形態によれば、重複レコードを検出するために、使用者に確認する回数を減らすことができる。

【0164】

また、使用者が、変換語候補が変換語であると判定した場合に増加した新たな重複レコードから変換語候補を抽出することができる。

【0165】

さらに、使用者が、重複レコード候補が重複レコードであると判定した場合に増加した新たな重複レコードから変換語候補を抽出することができる。

40

【0166】

なお、重複レコード判定部7は、重複レコードであると判定しなかった重複レコード候補が複数存在した場合、最も語の数の多い重複レコード候補から順に入出力部11に出力してもよい。すると、変換語候補抽出部10が変換語候補を抽出する可能性が高い順に重複レコード候補が入出力部11に出力されて使用者が重複レコードであるか否かを判定するため、変換語候補抽出部10が重複レコード候補から多くの変換語候補を抽出すると、語の数の少ない重複レコード候補からは変換語候補抽出部10が変換語候補を抽出する可能性が低くなり、変換語候補や、重複レコードを検出するために使用者に確認する回数を減らすことができる。

【0167】

また、変換語候補抽出部10は、使用者が複数の変換語候補を変換語であると判定する

50

と、使用者が判定した複数の変換語を類似度計算部 3 に出力し、類似度計算部 3 は、使用者が判定した複数の変換語に応じてデータベース 2 の該当するレコードを変換し、重複候補抽出部 6 は、変換されたレコードにもとづいて重複レコード候補の組を抽出してもよい。そして、重複レコード判定部 7 は、重複候補抽出部 6 が抽出した重複レコード候補の組のうち、重複レコード候補の組を構成する重複レコード候補の数が多い順に、重複レコード候補を入出力部 11 に出力してもよい。

【0168】

重複レコード判定部 7 は、使用者が複数の重複レコード候補を重複レコードであると判定すると、使用者が重複レコードであると判定した重複レコード候補を変換語候補抽出部 10 に出力し、変換語候補抽出部 10 は、抽出した変換語候補の数が多い重複レコード候補から抽出した変換語候補を、順に入出力部 11 に出力してもよい。

10

【0169】

実施の形態 5 .

本発明の第 5 の実施の形態を、図面を参照して説明する。図 13 は、本発明の第 5 の実施の形態の一構成例を示すブロック図である。

【0170】

本発明の第 5 の実施の形態の構成は、本発明の第 4 の実施の形態の構成に、重複レコードデータベース 9 に登録されている重複レコードの組を構成するレコードのうち、一のレコードを除いて、他のレコードをデータベース 2 から削除する重複レコード削除部 12 を加えたものであり、その他の構成要素は第 4 の実施の形態と同様のため、その他の構成要素には図 11 と同じ符号を付し、説明を省略する。

20

【0171】

重複レコード削除部 12 は、重複レコードデータベース 9 に登録された重複レコードの組を入出力部 11 に出力して、使用者に削除するレコードを選択させ、使用者が選択したレコードをデータベース 2 から削除してもよい。また、重複レコード削除部 12 は、重複レコードデータベース 9 に登録された重複レコードの組のうち、最もレコード ID の値の小さいレコード以外のレコードをデータベース 2 から削除してもよい。

【0172】

また、重複レコード削除部 12 は、削除したレコードの記録を記憶してもよい。

【0173】

以上に述べたように、この実施の形態によれば、重複レコード判定部 7 が、重複レコード判定ルールにもとづいて重複レコードであると判定した重複レコードを、重複レコードデータベース 9 に登録してから、重複レコード削除部 12 を介してデータベース 2 から削除するため、使用者が削除するレコードを確認したり、削除したレコードの記録を記憶させておいたりすることができる。

30

【0174】

実施の形態 6 .

本発明の第 6 の実施の形態を、図面を参照して説明する。図 14 は、本発明の第 6 の実施の形態の一構成例を示すブロック図である。

【0175】

本発明の第 6 の実施の形態の構成は、第 3 の実施の形態の構成に、使用者が新たにデータベース 2 に追加する情報を入力するデータベース登録部 13 と、データベース 2 に登録されている情報を検索する検索部 14 と、使用者に請求する料金を算出する検索料金算出部 15 とを加えたものであり、その他の構成要素は第 3 の実施の形態と同様のため、その他の構成要素には図 8 と同じ符号を付し、説明を省略する。

40

【0176】

データベース登録部 13 は、使用者が新たにデータベース 2 に追加する情報を入力すると、変換語辞書 5 に登録されている同義語と省略可能語ともとづいて、使用者が入力した情報と重複する情報である可能性のあるレコードを表示部に表示させる。

【0177】

50

例えば、図2の例に示す情報がデータベース2に登録され、図3の例に示す情報が変換語辞書5に登録されている場合に、使用者が、データベース登録部13に、名称が「日本電気奈良支店」である情報を入力する。

【0178】

すると、データベース登録部13は、入力された「日本電気奈良支店」に形態素解析等の方法を用いて、「日本電気奈良支店」を、「日本電気」と「奈良支店」との語に分解する。

【0179】

そして、データベース登録部13は、同義語辞書5を参照して、「日本電気」および「奈良支店」の同義語と省略可能語とを抽出する。「日本電気」の代表語である同義語は、「NEC」であるため、データベース登録部13は、「NEC」を抽出する。また、データベース登録部13は、「NEC」を代表語とする同義語である「エヌイーシー」を抽出する。

10

【0180】

さらに、データベース登録部13は、同義語辞書5を参照して、「日本電気」、「NEC」、および「エヌイーシー」のいずれかの省略可能語を抽出する。具体的には、「NEC」の省略可能語として「日電」を抽出する。

【0181】

そして、データベース登録部13は、同義語辞書5から抽出した語や、使用者が入力した情報の語を組み合わせて、使用者が入力した情報から変形した可能性がある情報を生成する。具体的には、「NEC奈良支店」、「エヌイーシー奈良支店」、「日電NEC奈良支店」、「日電エヌイーシー奈良支店」、「NEC日電奈良支店」、「エヌイーシー日電奈良支店」、「日電奈良支店」等を生成する。

20

【0182】

データベース登録部13は、生成した情報と合致する情報がデータベース2に登録されているか否かを検索して、データベース2から合致する情報を抽出する。すると、図2の例に示すレコードID「001」、レコードID「002」、およびレコードID「003」が抽出される。

【0183】

データベース登録部13は、抽出した各レコードを、重複可能性のあるレコードとして表示部に表示させる。

30

【0184】

検索部14は、上述したデータベース登録部13の動作と同様な動作を行って、データベース2から、使用者が検索部14に検索キーとして入力した情報、およびその情報から変形した可能性がある情報に合致する情報を、検索結果として表示部に表示させる。

【0185】

検索料金算出部15は、検索部14が検索結果を表示部に表示させると、使用者に請求する所定の料金を算出する。なお、使用者が、データベース2の所有者または管理者に、データベース2の使用料金を支払っている場合は、検索料金算出部15は、データベース2の所有者または管理者に請求する料金を算出してもよい。

40

【0186】

以上に述べたように、この実施の形態によれば、使用者がデータベース2に情報を登録する際に、重複する可能性のある情報を使用者に提示するため、新しく登録する情報が、重複レコードとなることを防ぐことができる。

【0187】

また、検索部14がデータベース2に登録されている情報を検索し、検索料金算出部15が、検索部14が行った情報の検索に応じた料金を算出するため、重複レコード検出システム20の所有者は、料金を使用者またはデータベース2の所有者または管理者に請求することができる。

【0188】

50

実施の形態 7 .

本発明の第 7 の実施の形態を、図面を参照して説明する。図 1 5 は、本発明の第 7 の実施の形態の一構成例を示すブロック図である。

【 0 1 8 9 】

本発明の第 7 の実施の形態の構成は、第 5 の実施の形態の構成に、重複レコード削除部 1 2 がデータベース 2 から削除した情報に応じて、データベース 2 の所有者等に請求する料金を算出する削除料金算出部 1 6 を加えたものであり、その他の構成要素は第 5 の実施の形態と同様のため、その他の構成要素には図 1 3 と同じ符号を付し、説明を省略する。

【 0 1 9 0 】

重複レコード削除部 1 2 は、第 5 の実施の形態における動作と同様の動作を行い、重複レコードをデータベース 2 から削除する。

【 0 1 9 1 】

削除料金算出部 1 6 は、重複レコード削除部 1 2 がデータベース 2 から削除した情報に応じて、データベース 2 の所有者または管理者に請求する料金を算出する。

【 0 1 9 2 】

以上に述べたように、この実施の形態によれば、データベース 2 に登録されている重複レコードの削除に応じた料金を、データベース 2 の所有者または管理者に請求することができる。

【 0 1 9 3 】

なお、重複レコード検出システム 2 0 の所有者等は、重複レコード検出システム 2 0 に他のデータベースに登録されている情報を入力して、重複レコード削除部 1 2 に情報を削除させてもよい。すると、変換語候補抽出部 1 0 が、変換語を変換語辞書 5 に登録するため、重複レコード判定部 7 による重複レコードの判定精度を向上させることができる。そのため、例えば、第 4 の実施の形態で、使用者が、重複レコード候補が重複レコードであるか否かの判定を行う回数を減らすことができる。

【産業上の利用可能性】

【 0 1 9 4 】

本発明は、データベースに重複して登録されている情報を抽出するシステムに適用することができる。

【図面の簡単な説明】

【 0 1 9 5 】

【図 1】本発明の第 1 の実施の形態の一構成例を説明するブロック図である。

【図 2】データベースに登録されている情報の例を示す説明図である。

【図 3】変換語辞書に登録されている情報の例を示す説明図である。

【図 4】本発明の第 1 の実施の形態の動作を説明するフローチャートである。

【図 5】本発明の第 2 の実施の形態の一構成例を示すブロック図である。

【図 6】重複判定ルールの例を示す説明図である。

【図 7】重複判定ルール記憶部が記憶している重複判定ルールの例を示す説明図である。

【図 8】本発明の第 3 の実施の形態の一構成例を示すブロック図である。

【図 9】本発明の第 3 の実施の形態の動作を説明するフローチャートである。

【図 1 0】重複レコード判定部が重複レコードであると判定したレコードの組の例を示す説明図である。

【図 1 1】本発明の第 4 の実施の形態の一構成例を示すブロック図である。

【図 1 2】本発明の第 4 の実施の形態の動作を説明するフローチャートである。

【図 1 3】本発明の第 5 の実施の形態の一構成例を示すブロック図である。

【図 1 4】本発明の第 6 の実施の形態の一構成例を示すブロック図である。

【図 1 5】本発明の第 7 の実施の形態の一構成例を示すブロック図である。

【符号の説明】

【 0 1 9 6 】

1 記憶部

10

20

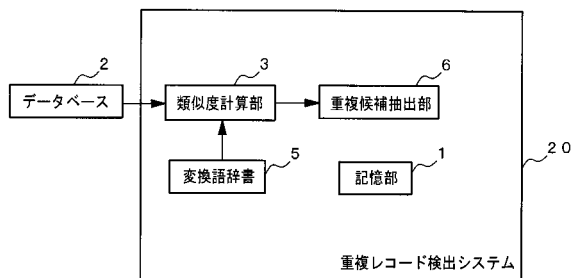
30

40

50

- 2 データベース
- 3 類似度計算部
- 4 データベース情報入力部
- 5 変換語辞書
- 6 重複候補抽出部
- 7 重複レコード判定部
- 8 重複判定ルール記憶部
- 9 重複レコードデータベース
- 10 変換語候補抽出部
- 11 入出力部
- 12 重複レコード削除部
- 13 データベース登録部
- 14 検索部
- 15 検索料金算出部
- 16 削除料金算出部
- 20 重複レコード検出システム

【図1】



【図2】

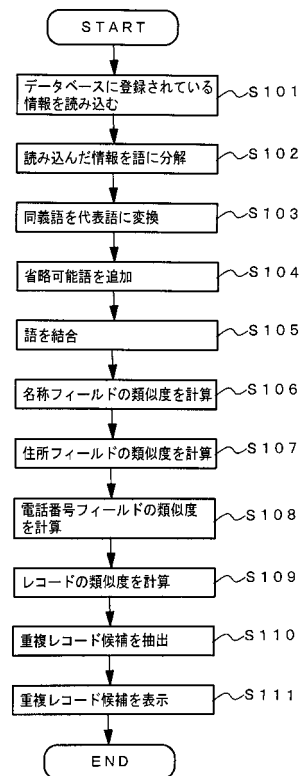
ID	名称	住所	電話番号
001	エヌイーシー奈良支店	〇〇〇1の1	000-111-1234
002	日電奈良支店	〇〇〇1-1	000-111-1235
003	NEC奈良支店	〇〇〇1-1	000-111-1234
...

【図3】

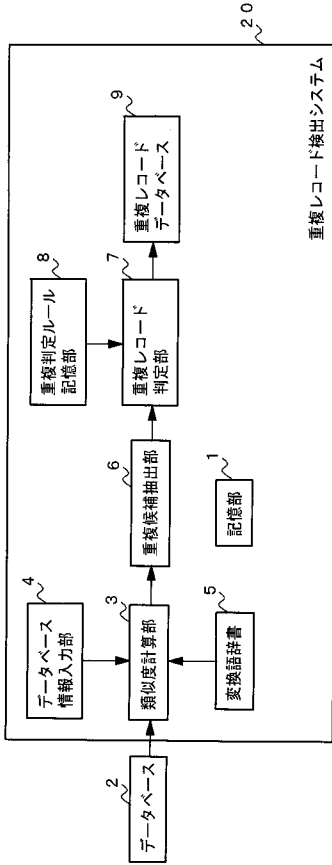
同義語辞書	
代表語	同義語
NEC	エヌイーシー
NEC	日本電気
日電	日本電気株式会社

省略可能語辞書	
省略可能語	グループID
NEC	1
日電	1

【図4】



【図5】



【図6】

重複判定ルール
レコード相互の類似度が特定の値を超えているならば、それらを重複レコードとみなす
レコード相互の類似度が特定の値以下であれば、それらを重複レコードではないとみなす
レコード相互の類似度が特定の値以下であれば、それらを重複レコードではないとみなす
いずれかのフィールドの類似度が特定の値以上であれば、それらを重複レコードではないとみなす
あるフィールドの類似度が特定の値以上であり、他のフィールドの類似度が所定の値以上であれば、
それらを重複レコードとみなす

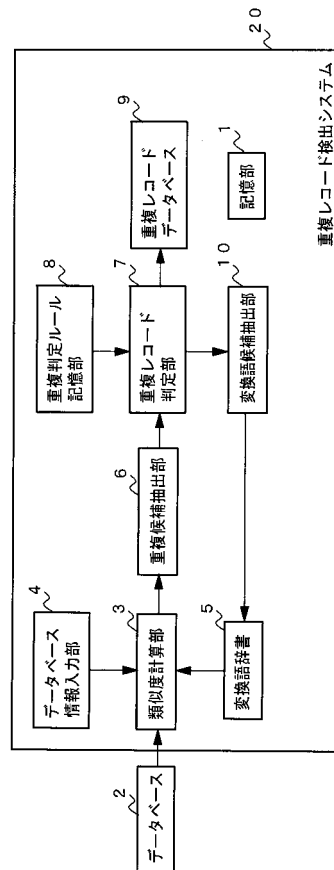
【図7】

```

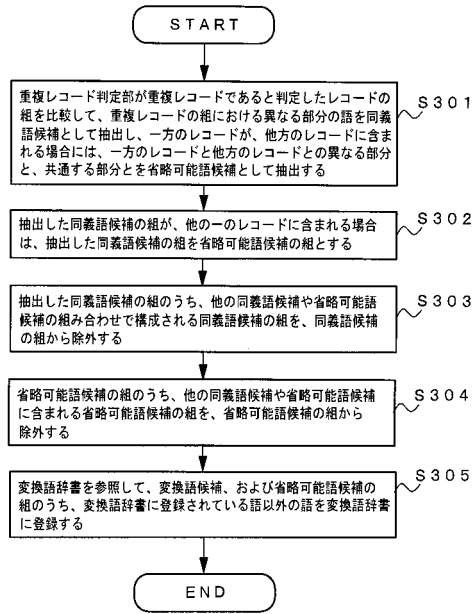
IF [類似度=1.0] THEN [重複レコードである] .....(1)
IF [ (住所フィールドの類似度>0.9)
AND (電話番号フィールドの類似度>0.9) ]
THEN [
IF [名称フィールドの類似度>0.9]
THEN [重複レコードである]
IF [名称フィールドの類似度<=0.9]
THEN [重複レコードではない]
] .....(2)

```

【図8】



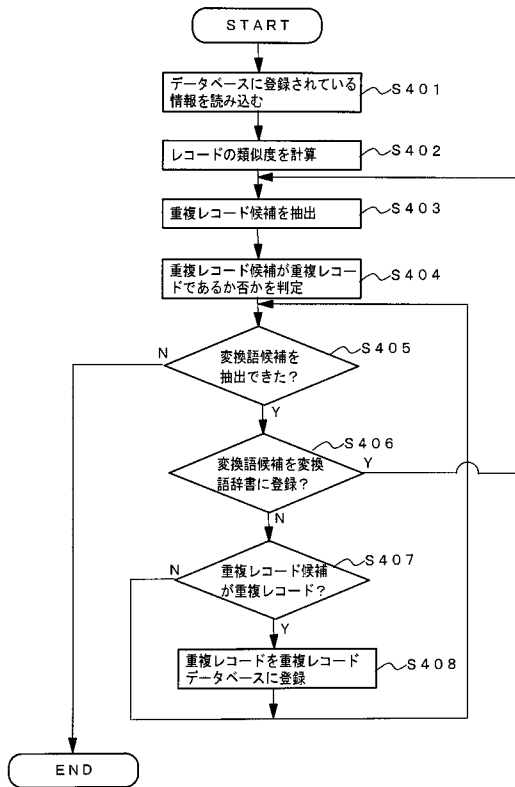
【図9】



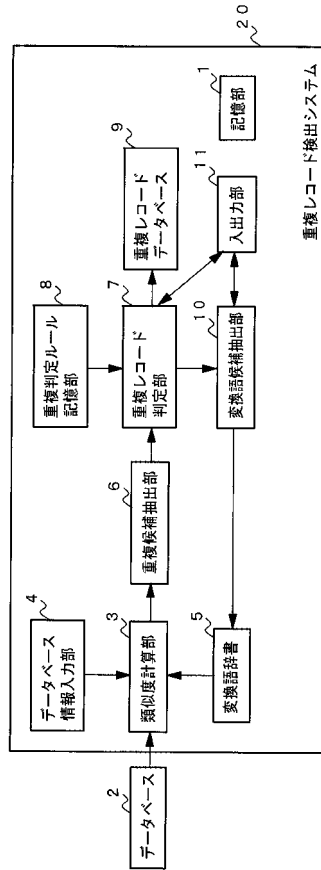
【図10】

ID	名称	住所	電話番号
001	日電NEC奈良支店	〇〇〇1-1	000-111-1234
002	エヌイーシー奈良支店	〇〇〇1-1	000-111-1234
003	日電エヌイーシー奈良支店	〇〇〇1-1	000-111-1234
004	NEC奈良支店	〇〇〇1-1	000-111-1234
005	日電奈良支店	〇〇〇1-1	000-111-1234

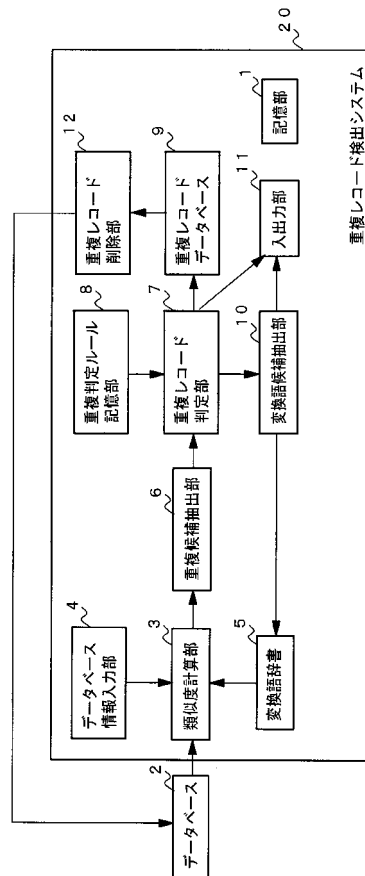
【図12】



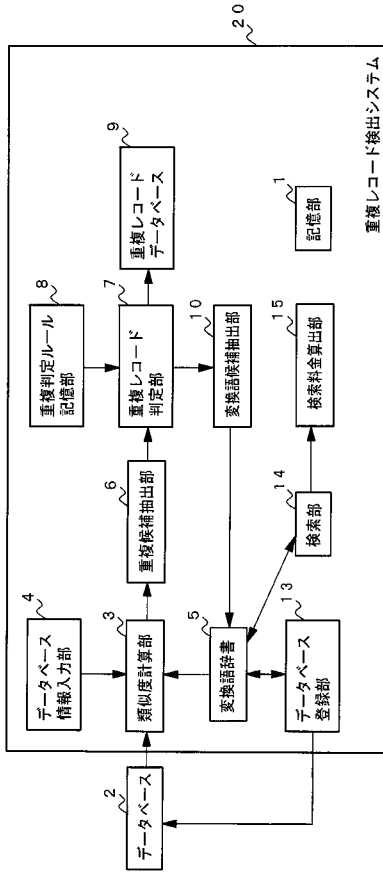
【図11】



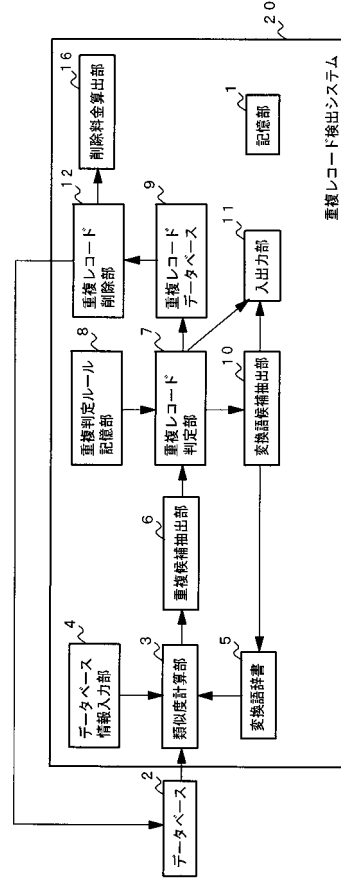
【図13】



【図14】



【図15】



フロントページの続き

(72)発明者 齋藤 悠
東京都港区芝五丁目7番1号 日本電気株式会社内

審査官 吉田 誠

(56)参考文献 特開平11-184884(JP,A)
特開平10-275159(JP,A)
特開平07-192053(JP,A)
特開平06-266769(JP,A)
特開2003-173345(JP,A)

(58)調査した分野(Int.Cl., DB名)
G06F 17/30