

(19) 日本国特許庁(JP)

(12) 特許公報(B2)

(11) 特許番号

特許第5535230号
(P5535230)

(45) 発行日 平成26年7月2日(2014.7.2)

(24) 登録日 平成26年5月9日(2014.5.9)

(51) Int. Cl.	F I				
G06N 3/00	(2006.01)	G06N	3/00	560A	
G06N 7/02	(2006.01)	G06N	7/02	554D	
G06F 17/30	(2006.01)	G06F	17/30	210D	

請求項の数 43 (全 54 頁)

(21) 出願番号	特願2011-533380 (P2011-533380)	(73) 特許権者	509123208
(86) (22) 出願日	平成21年10月23日(2009.10.23)		アビニシオ テクノロジー エルエルシー
(65) 公表番号	特表2012-507085 (P2012-507085A)		アメリカ合衆国 02421 マサチュー
(43) 公表日	平成24年3月22日(2012.3.22)		セッツ州 レキシントン スプリング ス
(86) 国際出願番号	PCT/US2009/061899		トリート 201
(87) 国際公開番号	W02010/048538	(74) 代理人	100079108
(87) 国際公開日	平成22年4月29日(2010.4.29)		弁理士 稲葉 良幸
審査請求日	平成24年10月23日(2012.10.23)	(74) 代理人	100109346
(31) 優先権主張番号	61/107,971		弁理士 大貫 敏史
(32) 優先日	平成20年10月23日(2008.10.23)	(72) 発明者	アンダーソン アーレン
(33) 優先権主張国	米国 (US)		イギリス オーエックスファイブ ツーエ
			スジー キドリントン イスリップ ロウ
			ワー ストリート 3
		審査官	塚田 肇
			最終頁に続く

(54) 【発明の名称】 ファジーなデータ操作

(57) 【特許請求の範囲】

【請求項1】

データ記憶システムに格納されたデータ要素をクラスタ化する方法であって、
上記データ記憶システムからデータ要素を読み取るステップ、
データ要素のクラスタを形成するステップであって、各データ要素が少なくとも1つのクラスタのメンバであり、且つ、所定のデータ要素における1つ以上のオブジェクトに少なくとも部分的に基づいて、上記所定のデータ要素が所定のクラスタのメンバとして割り付けられる、ステップ、

少なくとも1つのデータ要素を2つ以上のクラスタと関連付けるステップであって、上記2つ以上のクラスタのそれぞれに属する上記データ要素のメンバシップが、所定のクラスタにおける所定のデータ要素の内容と少なくとも1つの他のデータ要素の内容との間での比較に少なくとも部分的に基づいて、上記所定のクラスタにおける上記所定のデータ要素の部分的メンバシップを定量化する曖昧さの尺度によって表現され、上記所定のデータ要素の上記内容が、上記所定のデータ要素における1つ以上のオブジェクトにおける1つ以上の単語及び上記1つ以上の単語の変形の少なくとも1つを含む、ステップ、並びに

上記データ記憶システムに情報を格納して、上記形成されたクラスタを表現するステップ、

を含む方法。

【請求項2】

請求項1に記載の方法であって、

上記2つ以上のクラスタのそれぞれに属する上記データ要素のメンバシップを表現する上記曖昧さの尺度の各値が0と1との間にある、方法。

【請求項3】

請求項2に記載の方法であって、

上記メンバシップを表現する上記曖昧さの尺度の値が、上記2つ以上のクラスタのそれぞれに属する上記データ要素の確度に関連している、方法。

【請求項4】

請求項2に記載の方法であって、

上記2つ以上のクラスタのそれぞれに属する上記データ要素のメンバシップを表現する上記曖昧さの尺度の各値の合計が1である、方法。

10

【請求項5】

請求項4に記載の方法であって、

上記曖昧さの尺度の値を使用して会計の完全性を保つステップ、
を更に含む方法。

【請求項6】

請求項5に記載の方法であって、

所定の量についての会計の完全性を保つステップが、上記曖昧さの尺度の値によって上記量を重み付けするステップによって達成される、方法。

【請求項7】

請求項6に記載の方法であって、

上記メンバシップを表現する上記曖昧さの尺度の値を使用するデータ操作を実行するステップ、

20

を更に含む方法。

【請求項8】

請求項7に記載の方法であって、

上記データ操作が、上記1つ以上のクラスタの第1のクラスタ内の量の重み付けされた小計を計算するロールアップを含み、上記量が上記データ要素と関連付けられ、上記小計が、上記第1のクラスタ内で、上記第1のクラスタにおける上記データ要素の各々と関連付けられた上記量の値と上記第1のクラスタにおける上記データ要素の上記メンバシップを表現する上記曖昧さの尺度のそれぞれの値との積を合計することによって計算される、

30

方法。

【請求項9】

請求項8に記載の方法であって、

上記量の排他的小計及び上記量の包括的小計を計算するステップであって、上記排他的小計が、2つ以上のクラスタと関連付けられる上記第1のクラスタにおける上記データ要素を排除することによって計算され、上記包括的小計が、2つ以上のクラスタと関連付けられる上記第1のクラスタにおける上記データ要素を包含することによって計算されるステップ、

を更に含む方法。

【請求項10】

40

請求項2に記載の方法であって、

上記メンバシップを表現する上記曖昧さの尺度の値が関数に基づいて定められ、上記関数が、上記データ要素と上記2つ以上のクラスタとの間の関係を表現する、方法。

【請求項11】

請求項10に記載の方法であって、

上記関数によって表現される上記関係が、上記2つ以上のクラスタのそれぞれに属する上記データ要素の確度に関連している、方法。

【請求項12】

請求項10に記載の方法であって、

上記関数によって表現される上記関係が、上記データ要素と上記2つ以上のクラスタの

50

それぞれを表現する要素との間の定量化された類似性に基づく、方法。

【請求項 13】

請求項 12 に記載の方法であって、

上記 2 つ以上のクラスタのそれぞれを表現する上記要素が、それぞれのクラスタのキーである、方法。

【請求項 14】

請求項 1 に記載の方法であって、

上記 2 つ以上のクラスタの各クラスタに属する上記データ要素の上記曖昧さの尺度の値が、各クラスタについて等しい、方法。

【請求項 15】

請求項 1 に記載の方法であって、

上記 2 つ以上のクラスタの各クラスタに属する上記データ要素の上記曖昧さの尺度の値が、参照集合における上記データ要素の観測度数に基づく、方法。

【請求項 16】

請求項 1 に記載の方法であって、

上記 2 つ以上のクラスタの各クラスタが、上記データ要素における異なる潜在的エラーを表現し、上記 2 つ以上のクラスタの各クラスタに属する上記データ要素の上記曖昧さの尺度の値が、各クラスタによって表現される上記データ要素における上記潜在的エラーの確度に基づく、方法。

【請求項 17】

請求項 1 に記載の方法であって、

データ・クラスタを形成するステップが、
データ要素の複数のスーパークラスタを形成するステップ、及び、
各スーパークラスタについて、上記スーパークラスタ内にデータ要素のクラスタを形成するステップ、
を含む、方法。

【請求項 18】

請求項 17 に記載の方法であって、

各スーパークラスタを形成するステップが、
異なるデータ要素におけるオブジェクトの間での変形関係に基づいて、上記異なるデータ要素における上記オブジェクト間での整合を特定するステップ、
を含む、方法。

【請求項 19】

請求項 18 に記載の方法であって、

第 1 のオブジェクトと第 2 のオブジェクトとの間の変形関係が、予め定められた閾値未満にある上記第 1 のオブジェクトと上記第 2 のオブジェクトとの間の距離を表現する関数の値に対応する、方法。

【請求項 20】

請求項 19 に記載の方法であって、

上記変形関係が、同値関係ではない、方法。

【請求項 21】

請求項 17 に記載の方法であって、

少なくとも 1 つのデータ要素が、1 つより多くのスーパークラスタ中にある、方法。

【請求項 22】

データ記憶システムに格納されたデータ要素をクラスタ化するシステムであって、

上記データ記憶システムからデータ要素を読み取る手段、

データ要素のクラスタを形成する手段であって、各データ要素が少なくとも 1 つのクラスタのメンバであり、且つ、所定のデータ要素における 1 つ以上のオブジェクトに少なくとも部分的に基づいて、上記所定のデータ要素が所定のクラスタのメンバとして割り付けられる、手段、

少なくとも1つのデータ要素を2つ以上のクラスタと関連付ける手段であって、上記2つ以上のクラスタのそれぞれに属する上記データ要素のメンバシップが、所定のクラスタにおける所定のデータ要素の内容と少なくとも1つの他のデータ要素の内容との間での比較に少なくとも部分的に基づいて、上記所定のクラスタにおける上記所定のデータ要素の部分的メンバシップを定量化する曖昧さの尺度によって表現され、上記所定のデータ要素の上記内容が、上記所定のデータ要素における1つ以上のオブジェクトにおける1つ以上の単語及び上記1つ以上の単語の変形の少なくとも1つを含む、手段、及び

上記データ記憶システムに情報を格納して、上記形成されたクラスタを表現する手段、を備えるシステム。

【請求項23】

データ記憶システムに格納されたデータ要素をクラスタ化するためのコンピュータ・プログラムを格納するコンピュータ可読媒体であって、上記コンピュータ・プログラムが、上記データ記憶システムからデータ要素を読み取るステップ、

データ要素のクラスタを形成するステップであって、各データ要素が少なくとも1つのクラスタのメンバであり、且つ、所定のデータ要素における1つ以上のオブジェクトに少なくとも部分的に基づいて、上記所定のデータ要素が所定のクラスタのメンバとして割り付けられる、ステップ、

少なくとも1つのデータ要素を2つ以上のクラスタと関連付けるステップであって、上記2つ以上のクラスタのそれぞれに属する上記データ要素のメンバシップが、所定のクラスタにおける所定のデータ要素の内容と少なくとも1つの他のデータ要素の内容との間での比較に少なくとも部分的に基づいて、上記所定のクラスタにおける上記所定のデータ要素の部分的メンバシップを定量化する曖昧さの尺度によって表現され、上記所定のデータ要素の上記内容が、上記所定のデータ要素における1つ以上のオブジェクトにおける1つ以上の単語及び上記1つ以上の単語の変形の少なくとも1つを含む、ステップ、及び

上記データ記憶システムに情報を格納して、上記形成されたクラスタを表現するステップ、

をコンピュータに実行させるための命令を含む、コンピュータ可読媒体。

【請求項24】

請求項22に記載のシステムであって、

上記2つ以上のクラスタのそれぞれに属する上記データ要素のメンバシップを表現する上記曖昧さの尺度の各値が0と1との間にある、システム。

【請求項25】

請求項24に記載のシステムであって、

上記メンバシップを表現する上記曖昧さの尺度の値が、上記2つ以上のクラスタのそれぞれに属する上記データ要素の確度に関連している、システム。

【請求項26】

請求項24に記載のシステムであって、

上記2つ以上のクラスタのそれぞれに属する上記データ要素のメンバシップを表現する上記曖昧さの尺度の各値の合計が1である、システム。

【請求項27】

請求項26に記載のシステムであって、

上記曖昧さの尺度の値を使用して会計の完全性を保つ手段、を更に含むシステム。

【請求項28】

請求項27に記載のシステムであって、

所定の量についての会計の完全性を保つ手段が、上記曖昧さの尺度の値によって上記量を重み付けする手段によって達成される、システム。

【請求項29】

請求項28に記載のシステムであって、

上記メンバシップを表現する上記曖昧さの尺度の値を使用するデータ操作を実行する手

10

20

30

40

50

段、

を更に含むシステム。

【請求項 30】

請求項 29 に記載のシステムであって、

上記データ操作が、上記 1 つ以上のクラスタの第 1 のクラスタ内の量の重み付けされた小計を計算するロールアップを含み、上記量が上記データ要素と関連付けられ、上記小計が、上記第 1 のクラスタ内で、上記第 1 のクラスタにおける上記データ要素の各々と関連付けられた上記量の値と上記第 1 のクラスタにおける上記データ要素の上記メンバシップを表現する上記曖昧さの尺度のそれぞれの値との積を合計することによって計算される、システム。

10

【請求項 31】

請求項 30 に記載のシステムであって、

上記量の排他的小計及び上記量の包括的小計を計算する手段であって、上記排他的小計が、2 つ以上のクラスタと関連付けられる上記第 1 のクラスタにおける上記データ要素を排除することによって計算され、上記包括的小計が、2 つ以上のクラスタと関連付けられる上記第 1 のクラスタにおける上記データ要素を包含することによって計算される手段、を更に含むシステム。

【請求項 32】

請求項 24 に記載のシステムであって、

上記メンバシップを表現する上記曖昧さの尺度の値が関数に基づいて定められ、上記関数が、上記データ要素と上記 2 つ以上のクラスタとの間の関係を表現する、システム。

20

【請求項 33】

請求項 32 に記載のシステムであって、

上記関数によって表現される上記関係が、上記 2 つ以上のクラスタのそれぞれに属する上記データ要素の確度に関連している、システム。

【請求項 34】

請求項 32 に記載のシステムであって、

上記関数によって表現される上記関係が、上記データ要素と上記 2 つ以上のクラスタのそれぞれを表現する要素との間の定量化された類似性に基づく、システム。

【請求項 35】

請求項 34 に記載のシステムであって、

上記 2 つ以上のクラスタのそれぞれを表現する上記要素が、それぞれのクラスタのキーである、システム。

30

【請求項 36】

請求項 22 に記載のシステムであって、

上記 2 つ以上のクラスタの各クラスタに属する上記データ要素の上記曖昧さの尺度の値が、各クラスタについて等しい、システム。

【請求項 37】

請求項 22 に記載のシステムであって、

上記 2 つ以上のクラスタの各クラスタに属する上記データ要素の上記曖昧さの尺度の値が、参照集合における上記データ要素の観測度数に基づく、システム。

40

【請求項 38】

請求項 22 に記載のシステムであって、

上記 2 つ以上のクラスタの各クラスタが、上記データ要素における異なる潜在的エラーを表現し、上記 2 つ以上のクラスタの各クラスタに属する上記データ要素の上記曖昧さの尺度の値が、各クラスタによって表現される上記データ要素における上記潜在的エラーの確度に基づく、システム。

【請求項 39】

請求項 22 に記載のシステムであって、

データ要素のクラスタを形成する手段が、

50

データ要素の複数のスーパークラスタを形成する手段、及び、
各スーパークラスタについて、上記スーパークラスタ内にデータ要素のクラスタを形成する手段、

を含む、システム。

【請求項 40】

請求項 39 に記載のシステムであって、
各スーパークラスタを形成する手段が、

異なるデータ要素におけるオブジェクトの間での変形関係に基づいて、上記異なるデータ要素における上記オブジェクト間での整合を特定する手段、

を含む、システム。

10

【請求項 41】

請求項 40 に記載のシステムであって、

第 1 のオブジェクトと第 2 のオブジェクトとの間の変形関係が、予め定められた閾値未満にある上記第 1 のオブジェクトと上記第 2 のオブジェクトとの間の距離を表現する関数の値に対応する、システム。

【請求項 42】

請求項 41 に記載のシステムであって、

上記変形関係が、同値関係ではない、システム。

【請求項 43】

請求項 22 に記載のシステムであって、

少なくとも 1 つのデータ要素が、1 つより多くのスーパークラスタ中にある、システム

20

【発明の詳細な説明】

【技術分野】

【0001】

関連出願との相互参照

本願は、「ファジーなデータ操作 (FUZZY DATA OPERATIONS)」と題された米国特許出願第 61 / 107, 971 号 (2008 年 10 月 23 日出願。引用により本明細書に組み込まれる) に基づく優先権を主張する。

【0002】

30

本明細書は、データ管理の分野におけるファジーなデータ操作に関する。

【背景技術】

【0003】

クラスタ化、結合、検索、ロールアップ、及びソート等のデータ操作は、データ管理においてデータを取り扱うのに用いられる。クラスタ化は、種々のグループにデータを分類する操作である。結合は、2 つのデータを 1 つに結合する。キーによる検索は、そのキーに整合するデータ入力を見付け出す。ロールアップは、データのグループ全体の小計の 1 つ以上のレベル (又は他の組み合わせ) を計算する操作である。

【0004】

データ管理においてデータ品質は重要である。データ操作の結果として生ずる誤りや不正確がデータ品質を低下させる。例えば、ABC 社の従業員 (John Smith) を臨時雇用者として分類するか、又は常勤者として分類するかによって、異なるレベルの待遇が John Smith に与えられる。John Smith の雇用形態の誤った分類 (例えば、クラスタ化のデータ操作における誤り) は、ABC 社の人材データの品質に影響を及ぼす。

40

【0005】

データ操作の実装によっては、整合するレコードを識別したり、関連するレコードのグループを定義したり、又はレコードをリンクさせたりするのに、フィールド値 (「キー」) の厳密な比較を必要とする。データが曖昧であったり、不完全であったり、不十分であったり、又は不確実であったりする場合、フィールド値の厳密な比較に基づく方法は破綻

50

する場合がある。

【0006】

データ操作（例えば、クラスタ化）に関連する固有の曖昧さが存在する場合、当該曖昧さを解消するための1つのアプローチは、単に、当該曖昧さを無視することや、データを特定のグループに強制的に分類することであってもよい。例えば、ABC社の上記従業員（John Smith）は、マーケティング部門及びR&D部門の両方のために働いているとする。ABC社の人材データベースにおいては、John Smithは、マーケティング部門又はR&D部門の何れかと関連付けられているかもしれないが、多くの場合、一方の部門だけと関連付けられている。そのデータを特定のグループに強制的に分類することにより、当該固有の曖昧さを覆い隠して、データ品質に悪影響を及ぼす場合がある。

10

【0007】

ある事象（例えば、ある資産の所有権に関する事業体Aと事業体Bとの間での法的な争い）の結果が保留されているために、データ操作（例えば、クラスタ化）に関連する不確かさが存在する場合、データを特定のグループに強制的に分類することは、当該状況の流動性を扱う最善のアプローチではない場合がある。判決の前は、当該資産の所有権は不確かである。当該資産をA又はBの何れかに与えることは不正確であると後に判明するかもしれない。

【0008】

グループのメンバーシップの識別が曖昧であるために、データ操作（例えば、ロールアップ）に関連する不確かさが存在する場合、幾つかの選択肢の中の1つのグループにメンバーシップを割り付けて会計の完全性（accounting integrity）を保つことは、誤解を招く概念を生ずる場合がある。例えば、リスク評価や規制の目的のために取引先企業への貸出リスク（exposure on loans）を特定することに銀行が興味を持つ場合がある。取引先企業の識別は会社名によってなされる場合が多いが、このことは、会社名の記録形式に幅広いバラツキがあるために、曖昧な識別に繋がる場合がある。次に、このことは、取引先企業への貸出リスクの割り付けが曖昧であることを意味する。正しくは1つの会社に関連付けられる貸出が、実は当該1つの会社の名前の単に変形形式である、見掛け上別個の幾つかの会社にされるようなことが起こる場合がある。このことは、何れの単一の取引先企業に対しても当該銀行がリスク（exposure）

20

30

【0009】

正しくない情報又は欠落した情報のために、データ操作（例えば、結合）に関連する不確かさが存在する場合、データを特定のグループに強制的に分類したり、又は当該データを無視したりすることは、誤った関連付け又は情報の損失の何れかに繋がる場合がある。例えば、2つの異なるデータベースからのテーブルを結合しようとする際、これらのデータベースのテーブルによって共有される共通キーが全く無い場合が多い。このことを克服するために、これらのテーブル内のデータ（例えば、顧客の住所）を使用して、これら2つのデータベースにおけるレコード間の関係を推定する。しかしながら、住所の情報は正しくなかったり、又は不十分であったりする場合がある。ポータル・アドレス・ファイル（Postal Address File）等の信頼できる参照集合に照らした住所の妥当性確認により、テーブルAにおけるレコード上の番地（house number）が無効である（その番地を有する家は存在しない）一方、テーブルBには、その住所の有効な代替の完成形であるかもしれない複数の住所があると示されたと仮定する。テーブルAにおけるレコードにおける住所の完成形を任意に選択することは誤った関連付けに繋がる場合がある一方、当該レコードを無視することは情報の損失に繋がる。

40

【0010】

不正確なデータ入力のために、データ操作（例えば、検索）に関連する不確かさが存在

50

する場合、1つのアプローチは、単一の選択肢又は代替の補正值の単純なリストを提案することである。このことが、オペレータによってデータベースに入力されているデータの妥当性確認プロセスの一部である場合、複数の選択肢が存在する際に単一の選択肢を提案することは、その補正值を許容することにおいて誤った安心感をオペレータに与えることに繋がる場合がある。選択肢の単純なリストを提供する場合、これらの選択肢の中から選択するための合理的根拠をオペレータが有していない場合がある。単一選択が求められ且つ間違った選択に起因するデータ品質のある程度の低下が許容される場合、起こり得るデータ品質の損失の最小化及び定量化が目的となる。

【発明の概要】

【0011】

一般に、1つの態様において、データ記憶システムに格納されたデータ要素をクラスタ化する方法は、上記データ記憶システムからデータ要素を読み取るステップを含む。データ要素のクラスタが形成され、各データ要素は少なくとも1つのクラスタのメンバとなる。少なくとも1つのデータ要素が、2つ以上のクラスタと関連付けられる。上記2つ以上のクラスタのそれぞれに属するデータ要素のメンバシップは、曖昧さの尺度によって表現される。情報が上記データ記憶システムに格納され、形成されたクラスタを表現する。

10

【0012】

以下の特徴の1つ以上を含む態様もあり得る。

【0013】

上記2つ以上のクラスタのそれぞれに属するデータ要素のメンバシップを表現する曖昧さの尺度の各値がゼロ(0)と1との間にあってもよい。

20

【0014】

上記メンバシップを表現する曖昧さの尺度の値が、上記2つ以上のクラスタのそれぞれに属するデータ要素の確度に関連していてもよい。

【0015】

上記2つ以上のクラスタのそれぞれに属するデータ要素のメンバシップを表現する曖昧さの尺度の各値の合計が1であってもよい。

【0016】

上記方法は、上記曖昧さの尺度の値を使用して会計の完全性を保つステップ、を含んでもよい。

30

【0017】

所定の量についての会計の完全性が、上記曖昧さの尺度の値によって上記量を重み付けするステップによって達成されてもよい。

【0018】

上記方法は、上記メンバシップを表現する曖昧さの尺度の値を使用するデータ操作を実行するステップ、を含んでもよい。

【0019】

上記データ操作が、上記1つ以上のクラスタの第1のクラスタ内の量の重み付けされた小計を計算するロールアップを含み、上記量が上記データ要素と関連付けられ、そして上記小計が、上記第1のクラスタ内で、上記第1のクラスタにおけるデータ要素の各々と関連付けられた量の値と上記第1のクラスタにおけるデータ要素のメンバシップを表現する曖昧さの尺度のそれぞれの値との積を合計することによって計算されてもよい。

40

【0020】

上記方法は、上記量の排他的小計及び上記量の包括的小計を計算するステップであって、上記排他的小計が、2つ以上のクラスタと関連付けられる第1のクラスタにおけるデータ要素を排除することによって計算され、上記包括的小計が、2つ以上のクラスタと関連付けられる第1のクラスタにおけるデータ要素を包含することによって計算されるステップ、を含んでもよい。

【0021】

上記メンバシップを表現する曖昧さの尺度の値が関数に基づいて定められ、上記関数が

50

、上記データ要素と上記2つ以上のクラスタとの間の関係を表現してもよい。

【0022】

上記関数によって表現される関係が、上記2つ以上のクラスタのそれぞれに属するデータ要素の確度に関連していてもよい。

【0023】

上記関数によって表現される関係が、上記データ要素と上記2つ以上のクラスタのそれぞれを表現するデータ要素との間の定量化された類似性に基づく、方法。

【0024】

上記2つ以上のクラスタのそれぞれを表現する要素が、それぞれのクラスタのキーであってもよい。

【0025】

アレンジによっては、上記2つ以上のクラスタの各クラスタに属するデータ要素の曖昧さの尺度の値が、各クラスタについて等しくてもよい。

【0026】

上記2つ以上のクラスタの各クラスタに属するデータ要素の曖昧さの尺度の値が、参照集合におけるデータ要素の観測度数に基づいてもよい。

【0027】

上記2つ以上のクラスタの各クラスタが、上記データ要素における異なる潜在的エラーを表現してもよい。上記2つ以上のクラスタの各クラスタに属するデータ要素の曖昧さの尺度の値が、各クラスタによって表現されるデータ要素における潜在的エラーの確度に基づいてもよい。

【0028】

データ・クラスタを形成するステップが、
データ要素の複数のスーパークラスタを形成するステップ、及び、
各スーパークラスタについて、上記スーパークラスタ内にデータ要素のクラスタを形成するステップ、
を含んでもよい。

【0029】

各スーパークラスタを形成するステップが、
異なるデータ要素におけるオブジェクトの間での変形関係 (variant relation) に基づいて、上記異なるデータ要素におけるオブジェクト間での整合を特定するステップ、
を含んでもよい。

【0030】

第1のオブジェクトと第2のオブジェクトとの間の変形関係が、予め定められた閾値未満にある上記第1のオブジェクトと上記第2のオブジェクトとの間の距離を表現する関数の値に対応してもよい。

【0031】

アレンジによっては、上記変形関係は、同値関係 (equivalence relation) ではなくてもよい。

【0032】

少なくとも1つのデータ要素が、1つ以上のスーパークラスタ中にあるてもよい。

【0033】

もう1つの態様において、一般に、データ記憶システムに格納されたデータ要素をクラスタ化するシステムは、

上記データ記憶システムからデータ要素を読み取る手段、

データ要素のクラスタを形成する手段であって、各データ要素が少なくとも1つのクラスタのメンバである、手段、

少なくとも1つのデータ要素を2つ以上のクラスタと関連付ける手段であって、上記2つ以上のクラスタのそれぞれに属するデータ要素のメンバシップが、曖昧さの尺度によ

10

20

30

40

50

て表現される、手段、及び

上記データ記憶システムに情報を格納して、上記形成されたクラスタを表現する手段、を備える。

【0034】

もう一つの態様において、一般に、データ記憶システムに格納されたデータ要素をクラスタ化するためのコンピュータ・プログラムを格納するコンピュータ可読媒体について説明する。上記コンピュータ・プログラムは、

上記データ記憶システムからデータ要素を読み取るステップ、

データ要素のクラスタを形成するステップであって、各データ要素が少なくとも一つのクラスタのメンバである、ステップ、

少なくとも一つのデータ要素を二つ以上のクラスタと関連付けるステップであって、上記二つ以上のクラスタのそれぞれに属するデータ要素のメンバシップが、曖昧さの尺度によって表現される、ステップ、及び

上記データ記憶システムに情報を格納して、上記形成されたクラスタを表現するステップ、をコンピュータに実行させるための命令を含む。

【0035】

もう一つの態様において、一般に、キーを受け取り、データ記憶システムから一つ以上のデータ要素を返すデータ操作を実行する方法は、上記キーと上記データ要素の一つ以上の検索フィールドの値との間での候補整合(candidate matches)に基づいて、複数の候補データ要素を特定するステップを含む。上記候補整合は、上記検索フィールドとは異なる上記候補データ要素の一つ以上の比較フィールドの値に基づいて確認される。

【0036】

以下の特徴の一つ以上を含む態様もあり得る。

【0037】

上記データ操作が、

データ要素のクラスタを形成するステップであって、各データ要素が少なくとも一つのクラスタのメンバである、ステップ、を含んでもよい。

【0038】

少なくとも一つのデータ要素が二つ以上のクラスタと関連付けられ、上記二つ以上のクラスタのそれぞれに属するデータ要素のメンバシップが、曖昧さの尺度によって表現されてもよい。

【0039】

上記データ操作が、上記一つ以上のクラスタの第1のクラスタ内の量の重み付けされた小計を計算するロールアップを含み、上記量が上記データ要素と関連付けられ、そして上記小計が、上記第1のクラスタ内で、上記第1のクラスタにおけるデータ要素の各々と関連付けられた量の値と上記第1のクラスタにおけるデータ要素のメンバシップを表現する曖昧さの尺度のそれぞれの値との積を合計することによって計算されてもよい。

【0040】

また、上記方法は、上記量の排他的小計及び上記量の包括的小計を計算するステップであって、上記排他的小計が、二つ以上のクラスタと関連付けられる第1のクラスタにおけるデータ要素を排除することによって計算され、上記包括的小計が、二つ以上のクラスタと関連付けられる第1のクラスタにおけるデータ要素を包含することによって計算されるステップ、を含んでもよい。

【0041】

上記二つ以上のクラスタのそれぞれに属するデータ要素のメンバシップを表現する曖昧さの尺度の各値がゼロ(0)と1との間にあってもよい。

【0042】

10

20

30

40

50

上記メンバシップを表現する曖昧さの尺度の値が、上記2つ以上のクラスタのそれぞれに属するデータ要素の確度に関連していてもよい。

【0043】

上記メンバシップを表現する曖昧さの尺度の値が関数に基づいて定められ、上記関数が、上記データ要素と上記2つ以上のクラスタとの間の関係を表現してもよい。

【0044】

上記関数によって表現される関係が、上記2つ以上のクラスタのそれぞれに属するデータ要素の確度に関連していてもよい。

【0045】

上記方法は、所定のクラスタにおける所定のデータ要素の1つ以上の比較フィールドの値に基づいて、上記所定のデータ要素のメンバシップを特定するステップ、を含んでもよい。

10

【0046】

もう1つの態様において、一般に、キーを受け取り、データ記憶システムから1つ以上のデータ要素を返すデータ操作を実行するシステムは、

上記キーと上記データ要素の1つ以上の検索フィールドの値との間での候補整合に基づいて、複数の候補データ要素を特定する手段、及び

上記検索フィールドとは異なる上記候補データ要素の1つ以上の比較フィールドの値に基づいて、上記候補整合を確認する手段、
を備える。

20

【0047】

もう1つの態様において、一般に、キーを受け取り、データ記憶システムから1つ以上のデータ要素を返すデータ操作を実行するためのコンピュータプログラムを格納するコンピュータ可読媒体について説明する。上記コンピュータ・プログラムは、

上記キーと上記データ要素の1つ以上の検索フィールドの値との間での候補整合に基づいて、複数の候補データ要素を特定するステップ、及び

上記検索フィールドとは異なる上記候補データ要素の1つ以上の比較フィールドの値に基づいて、上記候補整合を確認するステップ、
をコンピュータに実行させるための命令を含む。

【0048】

もう1つの態様において、一般に、データ記憶システムにおけるデータ要素のデータ品質を測定する方法は、上記データ記憶システムからデータ要素(120)を読み取るステップ、を含む。上記データ要素の1つ以上のフィールドにおける1つ以上の入力の方々について、上記入力についての曖昧さの尺度の値が演算される。上記曖昧さの尺度の値に基づいて、上記データ記憶システムにおけるデータ要素のデータ品質の表現が出力される。

30

【0049】

以下の特徴の1つ以上を含む態様もあり得る。

【0050】

上記曖昧さの尺度の値を演算するステップが、上記データ要素の1つ以上のフィールドにおける入力を参照値と比較するステップ、を含んでもよい。参照値との厳密な整合ではない少なくとも第1の入力について1つ以上の変形が識別されてもよい。上記第1の入力についての変形に基づいて、上記第1の入力についての曖昧さの尺度の値が演算されてもよい。

40

【0051】

上記第1の入力についての曖昧さの尺度の値が、上記第1の入力についての変形の数に基づいてもよい。

【0052】

上記データ記憶システムにおけるデータ要素のデータ品質の表現が、特定数の変形を有する入力の数のヒストグラム・プロットを含んでもよい。

【0053】

50

上記特定数の変形が、ある範囲内にあるものとして規定されてもよい。

【0054】

上記データ記憶システムにおけるデータ要素のデータ品質の表現が、予め定められた閾値よりも大きい数の変形を有する入力のリストを含んでもよい。

【0055】

曖昧さの尺度の値を演算するステップが、1つ以上のフィールドにおける異なる入力のそれぞれの頻度を特定するステップ、を含んでもよい。他の入力の頻度と比較した第1の入力の相対頻度に基づいて、上記第1の入力についての曖昧さの尺度の値が演算されてもよい。

【0056】

もう1つの態様において、一般に、データ記憶システムにおけるデータ要素のデータ品質を測定するシステムは、

上記データ記憶システムからデータ要素を読み取る手段、

上記データ要素の1つ以上のフィールドにおける1つ以上の入力の各々について、上記入力についての曖昧さの尺度の値を演算する手段、及び

上記曖昧さの尺度の値に基づいて、上記データ記憶システムにおけるデータ要素のデータ品質の表現を出力する手段、
を備える。

【0057】

もう1つの態様において、一般に、データ記憶システムにおけるデータ要素のデータ品質を測定するためのコンピュータ・プログラムを格納するコンピュータ可読媒体について説明する。上記コンピュータ・プログラムは、

上記データ記憶システムからデータ要素を読み取るステップ、

上記データ要素の1つ以上のフィールドにおける1つ以上の入力の各々について、上記入力についての曖昧さの尺度の値を演算するステップ、及び

上記曖昧さの尺度の値に基づいて、上記データ記憶システムにおけるデータ要素のデータ品質の表現を出力するステップ、
をコンピュータに実行させるための命令を含む。

【0058】

もう1つの態様において、一般に、少なくとも1つのデータ記憶システムに格納された2つ以上のデータセットからのデータ要素を結合する方法は、第1のデータセットからのデータ要素におけるオブジェクトと第2のデータセットからのデータ要素におけるオブジェクトとの間の変形関係に基づいて、上記第1のデータセットからのデータ要素におけるオブジェクトと上記第2のデータセットからのデータ要素におけるオブジェクトとの間の整合を特定するステップ、を含む。それぞれのオブジェクトが整合として特定されたそれぞれのデータ要素が評価される。データ要素の上記評価に基づいて、上記第1のデータセットからのデータ要素が上記第2のデータセットからのデータ要素と結合される。

【0059】

以下の特徴の1つ以上を含む態様もあり得る。

【0060】

第1のオブジェクトと第2のオブジェクトとの間の変形関係が、予め定められた閾値未満である上記第1のオブジェクトと上記第2のオブジェクトとの間の距離を表現する関数の値に対応してもよい。

【0061】

上記変形関係は、同値関係ではなくてもよい。

【0062】

上記第1のデータセットからの第1のデータ要素におけるオブジェクトと上記第2のデータセットにおける第2のデータ要素におけるオブジェクトとの間の整合を特定するステップが、

上記変形関係が上記第1のデータ要素におけるオブジェクトと上記第2のデータ要素

10

20

30

40

50

におけるオブジェクトとの間に当てはまることを特定するステップ、
を含んでもよい。

【0063】

上記第1のデータセットからの第1のデータ要素におけるオブジェクトと上記第2のデータセットにおける第2のデータ要素におけるオブジェクトとの間の整合を特定するステップが、

上記変形関係が上記第1のデータ要素におけるオブジェクトと上記第1のデータセットにおける第3のデータ要素におけるオブジェクトとの間に当てはまること、及び上記変形関係が上記第3のデータ要素におけるオブジェクトと上記第2のデータ要素におけるオブジェクトとの間に当てはまること、を特定するステップ、
を含んでもよい。

10

【0064】

それぞれのオブジェクトが整合として特定された、それぞれのデータ要素を評価するステップが、

上記それぞれのデータ要素において整合として特定されたそれぞれのオブジェクト以外のオブジェクトを比較するステップ、
を含んでもよい。

【0065】

もう1つの態様において、一般に、少なくとも1つのデータ記憶システムに格納された2つ以上のデータセットからのデータ要素を結合するシステムは、

20

第1のデータセットからのデータ要素におけるオブジェクトと第2のデータセットからのデータ要素におけるオブジェクトとの間の変形関係に基づいて、上記第1のデータセットからのデータ要素におけるオブジェクトと上記第2のデータセットからのデータ要素におけるオブジェクトとの間の整合を特定する手段、

それぞれのオブジェクトが整合として特定されたそれぞれのデータ要素を評価する手段、及び

データ要素の上記評価に基づいて、上記第1のデータセットからのデータ要素を上記第2のデータセットからのデータ要素と結合する手段、
を備える。

【0066】

30

もう1つの態様において、一般に、少なくとも1つのデータ記憶システムに格納された2つ以上のデータセットからのデータ要素を結合するためのコンピュータ・プログラムを格納するコンピュータ可読媒体について説明する。上記コンピュータ・プログラムは、

第1のデータセットからのデータ要素におけるオブジェクトと第2のデータセットからのデータ要素におけるオブジェクトとの間の変形関係に基づいて、上記第1のデータセットからのデータ要素におけるオブジェクトと上記第2のデータセットからのデータ要素におけるオブジェクトとの間の整合を特定するステップ、

それぞれのオブジェクトが整合として特定されたそれぞれのデータ要素を評価するステップ、及び

データ要素の上記評価に基づいて、上記第1のデータセットからのデータ要素を上記第2のデータセットからのデータ要素と結合するステップ、
をコンピュータに実行させるための命令を含む。

40

【図面の簡単な説明】

【0067】

【図1】グラフに基づく演算を実行するシステムのブロック図である。

【図2A】複数のクラスタに属するデータ要素の例である。

【図2B】クラスタに対して実行される操作の例である。

【図2C】距離の計算例である。

【図2D】距離の計算例である。

【図3】ファジーなクラスタの図解である。

50

【図4】ファジーなクラスタのもう1つの図解である。

【図5】ファジーなクラスタを生成する方法のフローチャートである。

【図6】ファジーな検索の例の図解である。

【発明を実施するための形態】

【0068】

ファジーなデータ操作を実行するための手法は、データセットを格納する種々の形式のデータベースを含む様々なタイプのシステムに適用することができる。本明細書において使用されているように、データセットは、それぞれのフィールド（「属性」又は「カラム」とも呼ばれる）についての値を有する記録としてデータの部分を編成することを可能とする何等かのデータの集まりを含む。上記データベース・システム及び格納されたデータセットは、最新式のデータベース管理システム又は単純な単層ファイル（flat file）を格納するファイル・システム等の様々な形式の何れをとることもできる。様々なデータベース・システムの1つの態様は、データセット内のレコードに使用されるレコード構造のタイプである（各レコード内のフィールドに使用されるフィールド構造を含むこともできる）。システムによっては、データセットのレコード構造は、個々のテキスト文書をレコードとして単純に定義するものであってもよく、当該文書の内容は1つ以上のフィールドの値を表現する。システムによっては、単一のデータセット内の全てのレコードが同じ構造（例えば、フィールド構造）を有することを要件としない。

10

【0069】

複雑な演算は、有向グラフを通るデータ・フローとして表すことができることが多く（データフロー・グラフと呼ばれる）、当該演算の構成要素は当該グラフの節点（vertices）と関連付けられ、これらの構成要素の間のデータ・フローは当該グラフのリンク（アーク、エッジ）に対応する。このようなグラフに基づく演算を実装するシステムは、「グラフとして表された演算の実行（EXECUTING COMPUTATIONS EXPRESSED AS GRAPHS）」と題された米国特許第5,966,072号に記載されている（引用により、本明細書に組み込まれる）。グラフに基づく演算を実行するための1つのアプローチは、当該グラフの種々の節点と各々が関連付けられた多数のプロセスを実行すること、及び当該グラフのリンクに従って当該プロセス間の通信経路を確立することである。例えば、上記通信経路は、TCP/IP若しくはUNIX（登録商標）ドメイン・ソケットを使用するものであっても、プロセス間でデータを渡すのに共有メモリを使用するものであってもよい。

20

30

【0070】

図1を参照すると、グラフに基づく演算を実行するシステム10は、データ・ストア12に連結された開発環境13及びデータ・ストア12に連結されたランタイム環境18を備える。開発者11は、開発環境14を使用してアプリケーションを構築する。アプリケーションは、データ・ストア12におけるデータ構造（開発者による開発環境14の使用結果として当該データ・ストアに書き込まれたものであってもよい）によって規定される1つ以上のデータフロー・グラフと関連付けることができる。演算グラフ15についてのデータ構造13は、例えば、演算グラフの節点（構成要素又はデータセット）及び節点間のリンク（ワーク要素のフローを表現する）を規定する。また、上記データ構造は、データフロー・グラフの構成要素、データセット、及びフローの様々な特徴を含むこともできる。

40

【0071】

ランタイム環境18は、UNIX（登録商標）オペレーティング・システム等の、好適なオペレーティング・システムの制御下にある1つ以上の汎用コンピュータ上にホスティングすることができる。例えば、ランタイム環境18は、ローカル型であっても（例えば、SMPコンピュータ等のマルチプロセッサ・システム）、又はローカル分散型であっても（例えば、クラスタ若しくはMPPとして連結された複数のプロセッサ）、又はリモート型であっても、又はリモート分散型であっても（例えば、LAN若しくはWANネットワークを介して連結された複数のプロセッサ）、あるいはそれらの何れの組み合わせであ

50

っても、複数の中央処理装置（CPU）を使用するコンピュータ・システムの構成を備える多重ノード並列演算環境を含むことができる。

【0072】

ランタイム環境18は、データ・ストア12及び/又はユーザ17から演算を実行及び構成するための制御入力を受け取るように構成される。上記制御入力は、（格納されたグラフ・データ構造において規定される）対応するデータフロー・グラフを使用して特定のデータセットを処理するコマンドを含むことができる。ユーザ17は、例えば、コマンド・ライン又はグラフィカル・インターフェースを使用して、ランタイム環境18と双方向に働きかけることができる。

【0073】

ランタイム環境18は、前実行モジュール20及び実行モジュール22を含む。前実行モジュール20は、あらゆる前処理プロシージャを実行し、様々なファジーな操作（例えば、米国特許出願公開第2009/0182728号に記載されているような操作。引用により、本明細書に組み込まれる）に使用される辞書21やアーカイブ24等の演算グラフを実行するためのリソースを準備し、維持する。辞書21は、データセットに現れる単語及び単語についての関連情報を格納する。アーカイブ24は、データセットの単語、句、又はレコードに基づく前処理からの様々な結果を格納する。辞書21及びアーカイブ24は、様々なフォーマットの何れでも実装されてることができ、データの単一の集まりとして、又は複数の辞書及びアーカイブとして、編成することができる。実行モジュール22は、構成要素の演算を実行するために演算グラフに割り付けられたプロセスの実行をスケジューリングし、制御する。実行モジュール22は、グラフの構成要素と関連付けられる処理中にアクセスされるシステム10に連結された外部の演算リソース（例えば、データベース・システムからレコードを提供するデータ・ソース26）と双方向に働きかけることができる。システム10において実行されるファジーな操作は、データ品質を評定するためのデータの分析又はデータの編成及び/若しくは統合等の様々な目的に使用することができる。

【0074】

何れの企業又は他の組織においても、核となる資産は、取引する製品、サービス及び顧客、個人、銀行及び他の企業との契約及び会計のリストに由来する、その事業を運営するためにその企業又は他の組織が保有しているデータである。このデータは、複数のフォーマットで、紙やスプレッドシートからリレーショナル・データベースやエンタープライズ・アプリケーションに至る複数のシステム（例えば、会計システム又はサプライ・チェーン管理システム）に、格納される。全ての組織の懸念の中心は、このデータの品質及び完全性である。

【0075】

請求書が正しくない価格や偽装された商品を含む場合、間違った金額が課されたり、又は誤った品物が配送されたりする。顧客又はサプライヤの住所が間違っている場合、出荷又は発注が遅れたり又は無くなったり、請求書又は支払いが目的とする相手に届かなかったりする場合がある。あるシステム上での顧客を表現するキーが別のシステム上での異なる顧客の会計にリンクしている場合、その顧客の会計の状態についての報告は信頼できないであろうし、更に悪いことに、ある顧客が他の顧客の会計へのアクセスを有することになる場合もある。不十分なデータ品質は、事業のきちんとした運営を混乱させ、収益の損失、悪評又は機会損失に繋がる場合がある。

【0076】

事業又は組織のデータの重要な部分集合は、時としてマスター・データと称される、非取引（non-transactional）参照用データである。これは、製品、顧客、会計、サプライヤ、及び各データ項目の特定の属性を表現するのに使用される特定の有効な値（例えば、顧客が男性若しくは女性の性別を有すること、又は製品が列举リストの中の1つの色を有すること）のリストを含むことができる。一般に、取引や価格等の、組織の短期的な運営データはマスター・データから排除される。マスター・データ管理は、

10

20

30

40

50

組織及び組織の参照用データのメンテナンスと関係する。

【0077】

データ品質及び参照整合性に関する問題は多くの形式をとる。これらの問題は、一貫性を保つのが困難な場合がある、異なる種類のデータ・システムが複数存在することによって悪化する。潜在的な問題の非網羅的なリストは以下の通りである。

【0078】

1) データが不正確に入力又は記録される場合がある：なされた入力が、意図したものではない。例えば、顧客の名前若しくは住所、製品のラベル若しくは説明書、又は列挙リストから選ばれるべき値における単語のスペル違いに繋がる誤植又は転記の誤りが入力にある場合がある。多くのデータ入力アプリケーションは、入力時にデータを確認して、このようなエラーを防ぐことを目的とするセーフガードを有するが、未だにエラーは起きている。

10

【0079】

2) データが不十分である場合がある：全てのフィールドが埋まっていない。顧客アプリケーションが特定のフィールドが欠落している情報を有していたかもしれない。フォームの完成が入力中に中断されたかもしれない。情報が入力時に無効とみなされ、破棄されたかもしれない。入力を完成するための情報が入力時には入手可能ではなかったかもしれない(何等かの他の作業の完了の保留等)。

【0080】

3) データが無効である場合もある：フィールドは埋まっているが、無効な値で埋まっている。列挙リストから期待される値の何れにも入力が整合しない場合がある。期待されるデータ型に対して入力が有効ではないかもしれない(例えば、十進数のフィールドに英字が存在したり、日付における日が、その月の日数よりも大きかったり(例えば、6月31日)する場合がある)。

20

【0081】

4) 誤ったフィールドにデータが入力される場合がある：住所の通りのフィールドに、都市名コードや郵便番号が現れる場合がある。期待されるものと異なるフォーマットを有する外国の住所を、期待されるフォーマットに強制的に適合させたかもしれない。請求書又は発注書のフォーム上の説明又はコメントのフィールドに製品IDがあるかもしれない。名字が一般的な名である場合(例えば、Gregory Paul)や、名前が珍しい名前若しくは外国の名前である場合、個人の姓(ラスト・ネーム)及び名(ファースト・ネーム)が入れ替わっているかもしれない。

30

【0082】

5) データ入力についての標準が存在しない場合がある：データが一貫性無く入力される場合がある。住所の行の順序が標準化されておらず、同じデータセットにおいてさえ、常に同じように記録されている訳ではない場合がある。会社名の詳細なフォームが標準化されておらず、同じデータセットにおいてさえ、多くの変形形式が許容されている場合がある。顧客名が正式なミドル・ネームを含む場合があり、又はミドル・ネームが無い場合があり、又はミドル・イニシャルだけの場合もある。同様に、名(ファースト・ネーム)がイニシャルだけの場合もある。「二重姓」の姓(ラスト・ネーム)が、ハイフンを伴う若しくは伴わない姓(ラスト・ネーム)に存在する場合があり、又は、ミドル・ネームのフィールドと姓(ラスト・ネーム)のフィールドとの間で分けられている場合がある。

40

【0083】

6) データがフリー・テキスト・フィールドに保持されている場合がある：請求書又は発注書のフォーム上のコメント・フィールドにおける注釈が、製品名又は記述的な属性等の、さもなければ欠落していたであろう、重要な情報を含んでいる場合がある。例えば、結婚により女性の姓(ラスト・ネーム)が変わる場合、データベース・テーブルにおける記述フィールドが、他のフィールドに対する変更についての説明を含む場合がある。

【0084】

7) キー関係が壊れている場合がある：データベースはキーを使用して、異なるテーブ

50

ルや時には異なるデータベースにおけるテーブルに保持されている関連データをリンクする。キーが正しいレコードを正しくリンクしない場合、当該データベースの参照整合性が破綻している。ある顧客が正しくは別の顧客に属する会計にリンクされているような場合、キーがレコードを不正確にリンクする場合がある。存在しないレコードにキーがリンクする場合もある（例えば、ある会計レコードの顧客キーが、実在する顧客レコードの何れにもリンクしない）。場合によっては、上記顧客レコードは存在するが、異なるキーを有する（このようなキーは時として「紛失キー」と記述される）。他の場合においては、対応する顧客レコードが全く存在しない場合、当該会計レコードは孤児であると言われる。

【 0 0 8 5 】

8) キー関係が存在しない場合がある：出所が異なるデータベースが類似のデータを保持するものの、それらによって共有されるデータをリンクするキーを保持していないことがある。ある事業分野が別の事業分野との顧客の共有を実現しないことは希である。企業又は組織が合併する場合、これら2つの事業体のマスター・データが組み合わせられる場合がある。これら2つの事業体の異なる標準及び不等価な有効値により、マスター・データの一貫した集合の達成が困難となるが、顧客等の共有データの識別及びリンクの問題の方がより厳しいことが多い。

【 0 0 8 6 】

データ・クレンジングは、これらの問題点を識別し、訂正しようとするものである。レガシー・システムの数及び複雑さ、システム間のインターフェースの数、並びに新しいシステムの導入の速度の故に、真の課題は、データ品質の問題を如何に修復するかではなく、それらを如何に上手く対処するかであることが多い。

【 0 0 8 7 】

企業又は組織のシステムにおけるデータを見付け出し、アクセスし、そして操ることににおける主要な概念は、「キー」の概念である。主キーはフィールド、又はフィールドの組み合わせであり、その値は、データセットにおけるレコードを一意的に識別するのに役立つ。リレーショナル・データベース内では、各テーブルが、そのテーブル内のレコードを一意的に識別する主キーを有してもよい（主キーが一意的でない場合、それはデータ品質の別問題である）。テーブルにおける外部キーは、他のテーブルにおけるレコードにリンクするキーである。

【 0 0 8 8 】

多くのデータ操作を行うことができるが、データ操作は、データベースのテーブル又は他のデータセットのキーに依存する。よく知られているキーに基づくデータ操作は、キーによるルックアップ、結合、ロールアップ、スキャン、ソート、マージ及び、並列処理における、キーによる分割である。これらのデータ操作は、キーの厳密な一致（*agreement*）（本明細書においては「厳密な整合」と呼ぶ）に基づく。ルックアップというデータ操作においては、キーを使用して、厳密に整合するキーを有するルックアップ・データセットから1つ以上のレコードを読み出す。結合というデータ操作においては、1回に1レコードずつ、あるデータセットからのレコードの内容に、別のデータセットからの共通キーを共有するレコードの内容を付け加える（そして、ことによると、部分集合化することによって、2つ（又は2つ以上）のデータセットを組み合わせる。2つ以上のレコードが整合する共通キーを有する場合、整合するレコードのペアの各々について別個の出力レコードが形成される。

【 0 0 8 9 】

ロールアップというデータ操作においては、共通キーを共有するレコードのグループの内容を組み合わせ、同じキーを有する単一の出力レコードを生成する。例えば、取引金額を合計しながら取引レコードを顧客レベルにロールアップすることにより、ある顧客の全取引金額が得られる。スキャンというデータ操作においては、共通キーを共有するレコードのグループにおける各レコードについて、同じキーを有する、これまでに見られた全てのレコードの内容を使用して、出力レコードを算出する。例えば、顧客による取引のス

10

20

30

40

50

キャンを用いて、顧客支出の現在高を算出することができる。

【0090】

ソートというデータ操作においては、レコードをそれらのキー値によって順序付ける。マージというデータ操作においては、1つ以上のデータ・ストリームからのソートされたデータを組み合わせて、出力レコードもまたソートされているように、単一のストリームとする。並列処理におけるキーによる分割というデータ操作においては、キーの値に基づいて、データを処理パーティションに割り当てる。

【0091】

上記において考察した種類のデータ品質問題を各々が有するかもしれない複数の独立したシステムが共存する場合、レコードを共通データと関係付けるキーは一般的には存在せず、存在するキーは信頼性が無い場合がある。結局のところ、各レコード内のデータは重要な項目である。キーは、データを識別し、データにアクセスするためにデータベースに導入される便利なフィクションであると考えられる。信頼性のあるキーが無い場合、識別の目的にデータそのものを使用してもよい。

10

【0092】

レコードの内容によるレコードへのアクセスは検索に基づく。例えば、1つのデータベースにおける顧客を、第2のデータベースにおいて名前によって探してもよい。名前は曖昧な識別子であるので、それをキーとすることは希である。名前を使用して識別を開始することができるけれども、整合を確証するのに誕生日や住所等の支援情報が必要とされるのが一般的である。

20

【0093】

その上、データ品質問題の故に、レコードが正しく整合するために、名前も確証情報も厳密に一致することが必要とされない。厳密な一致は限定的すぎる場合があり、正確な整合を要求することが、多くの正しい識別の欠落に繋がる場合がある。(ファジーな)検索のデータ操作により、しっかりと整合するけれども、必ずしも厳密には整合しないデータ入力を取り出される。例えば、「Leslie」についてのファジーな検索により、「Lesley」と名付けられた人物についてのレコードが返される場合がある。ファジーな検索において、種々の度合いの類似性又は確証を有する整合レコードが2つ以上存在する場合がある(Laslieについての検索により、Lesleyと名付けられた第2の人物のレコードが取り出される場合もある)。候補整合が十分に確証されず、限定的整合又は許容整合と称される場合がある。例えば、取り出されたLesleyのレコードの誕生日がLeslieのレコードの誕生日と一致しない場合があり、この場合、当該候補整合は確証されない。

30

【0094】

検索時、ルックアップのための厳密なキーを使用する単一段階プロセスを、2段階プロセスによって置き換える。レコードは、検索語を使用する取り出しのために識別され、整合を特定するために評価される。検索語がレコードを一意的に識別することは希であるので、検索語はキーではない。しかしながら、検索語をキーのように使用して、レコードをリンクする。

【0095】

明確にするために、そこから検索語を選択するためのフィールドを、レコードを比較して整合の品質を評価するのに使用されるフィールドと区別することが有用である。これらを、それぞれ検索フィールド及び比較フィールドと呼んでもよい。

40

【0096】

検索語又は比較フィールドが同一ではない場合に整合を見出して判定するために、採点関数を使用して、変形値(variant value)を認識してもよい。候補整合を、変形検索語を使用して取り出し、採点関数を使用して評価して、確証フィールド間の整合の品質を定量化してもよい。これらの採点関数は、様々なデータ品質問題を明らかにするように設計される。それらは、これらの問題にもかかわらず、減点されたスコアで整合を認識する。例えば、個人名についての採点関数は、姓(ラスト・ネーム)と名(ファ

50

スト・ネーム)との入れ替え又はミドル・イニシャルの使用には寛容であるようにしてもよく、企業名用に調整された採点関数は、欠落している単語よりも単語の順序に重きを置くようにしてもよい。

【0097】

厳密なキーのもう1得tの基本的な用途は、共通するキー値を有するレコードの集合(キーグループと呼ばれることが多い)を識別することである。これらのキーグループは、多くのキーに基づくデータ操作において重要な役割を演ずる。厳密に整合するキーの要件が緩和される場合、如何にしてキーをグループ化するかという疑問が生ずる。緩和された整合基準に基づいて一緒にグループ化されたキーの集合は、クラスタと呼ばれる。

【0098】

一般に、クラスタは、それらの比較フィールドが比較テストに合格するレコードの集合であってもよい。例えば、1つのアレンジにおいて、クラスタについてのあるレコードのスコアが閾値を超える場合、そのレコードはそのクラスタのメンバである。クラスタについてのレコードのスコアを定義する方法としては種々の方法が多数存在し、概して、クラスタの各メンバについてレコードを個々に採点し、次いでそれらのスコアを足し合わせることを含むけれども、この方法に限定されるものではない。例えば、スコアは、クラスタの各メンバについてのレコードのスコアの最大値であってもよく、又はクラスタの各メンバについてのスコアの平均値であってもよい。アレンジによっては、レコードのペアを採点することは、フィールド値の一方の集合を他方の集合と比較した結果に数を割り付けることを含む。フィールド値の比較が、定性的評価及び定量的評価の両方を含んでもよい。

【0099】

クラスタの定義における問題は、フィールド値の比較が、採点された関係であるため、曖昧なメンバシップの割り当てが可能であるということである。特に、採点は、2つ以上のクラスタに属する1つのデータを示す場合がある。1つのアレンジにおいて、上記1つのデータを上記クラスタの中の1つに強制的に属させて、それらが厳密なキーの場合にあるとして、それらのクラスタをハッキリと定義することによって、この曖昧さに対処してもよい。この場合、上記キーに基づくデータ操作は、本質的に、それらが厳密なキーの場合にあるときのみである。

【0100】

厳密なキーに基づくデータ操作は、様々な理由により、必ずしも望まれるほどには正確又は精密ではない場合がある。1つの理由は、1つのデータ又はデータ操作に関連する固有の曖昧さである場合がある。例えば、1つのデータが、2つ以上のグループに正当に属する場合がある。クラスタ化の方法によっては、固有の曖昧さにより、精密な分類が困難又は達成不能となる場合がある。例えば、従業員が、彼又は彼女が属する部門に従って分水される、上述の人材データベースにおいて、従業員は、マーケティングとR&D等、同時に2つの部門に属する場合がある。当該従業員をいずれかの部門(マーケティング又はR&D)に強制的に関連付けることは、誤解を招くことになる場合がある。単純に、当該従業員を両方の部民に関連付けることは、二重カウント問題を引き起こす場合がある。例えば、医療費等の費用が、同じ従業員に2回記録される場合がある。

【0101】

正確な分類が可能ではない場合があるもう1つの理由は、保留中の事象の結果が現行の分類に強い影響を与える場合があることである。例えば、慈善団体か又は非慈善団体かという組織の法的地位により、その納税義務に変化が生ずる場合がある。IRS(国税庁)と当該組織との間に、当該組織が慈善団体として適格であるが故に税控除を受ける資格があるか否かに関する進行中の訴訟が存在すると更に仮定する。当該組織の年間予算において、当該組織の課税上の地位が慈善団体の地位にあると推定され、従って納税を保留して、より小さい予算が設定されている場合、且つ、その後、当該組織は、非慈善団体であるが故に慈善団体のみ与えられる税控除を受けることができないと裁判所が決定する場合、当該年間予算は改訂されなければならない。このような状況に対処する旧来の方法は、当該予算に影響を及ぼす可能性のある不都合な判決を説明する注記を予算に付加すること

10

20

30

40

50

であることが多い。不都合な判決が下された場合、当該予算は改定されなければならない。しかしながら、予算を修正しなければならないことよりも更に悪いことに、当該予算が、他の事業分野、又は他の国において、他のアプリケーションによって使用されていた場合は、その波及効果を追跡するのが不可能である場合があるため、当該予算そのものを超える修正が不可能である場合がある。

【0102】

上記2つの例は、従来のデータ・アプローチが、クラスタへの曖昧な割り付け（「部分的メンバシップ」）を取り扱うのに如何に妥当ではない場合があるかを説明している。曖昧な割り付けは、クラスタへの1対1整合を、確実にすることができないか又は確実にしない方がよい場合に起こる。複数のクラスタへの割り付けを許すことによって呈される1つの難題は、会計の完全性を如何に保つかということである。部分的メンバシップの方法は、この目的のために使用することができ、この開示において後に詳細に考察する。曖昧なメンバシップを有する要素を含むクラスタを取り扱うのに、ファジーなデータ操作を使用してもよい。

10

【0103】

重複するメンバシップを伴う要素をクラスタが有し且つ2つ以上のクラスタに何等かのデータが関連付けられている場合、ファジーなロールアップというデータ操作を使用して、会計の完全性を保ち且つ可能な代替の割り付けと関連するエラーの範囲を報告しつつ、計算を実行してもよい。従業員が2つ以上の部門のために働いている場合、その従業員のための費用を、当該従業員の部分的メンバシップに従って、それらの部門の間で配分することができる。

20

【0104】

クラスタ・メンバシップが、上記法的な例におけるように、将来の事象次第であるか、又は上記銀行業務における取引先企業の特定期間におけるように、曖昧又は不十分な情報のために不確実であるか、の何れかである場合、例えば、グループによる金銭的な合計を演算するファジーなロールアップ操作は、会計の完全性を保ちつつ、この不確実さを反映すべきである。確かに、会社の税分類等の、不確実な将来の事象の場合、ある偶発性が生ずる。特定の選択肢への早まった割り付けは、企画立案及びリスク評価の目的に誤解を生ずる概念を生ずる場合がある。

【0105】

例えば、図2Aにおいて、データ要素120がクラスタ122、クラスタ124、又はクラスタ126の何れに属するかは不確実である。データ要素120がこれら3つのクラスタ122、124、及び126に同時に属することもあるかもしれない。また、データ要素120が、ある特定の時刻には1つのクラスタに属するけれども、上記3つのクラスタの間で交代することもあるかもしれない。クラスタ122、クラスタ124、及びクラスタ126に属するデータ要素120のメンバシップは n_1 、 n_2 、及び n_3 によって表現され、 n_1 、 n_2 、及び n_3 は分数である。データ要素120が等しい可能性で同時に3つのクラスタに属する場合、 n_1 、 n_2 、及び n_3 の各々には分数 $1/3$ を割り付けることができる。この場合、クラスタ122、124、及び126に属するデータ要素120の部分的メンバシップの合計は $1(1/3 + 1/3 + 1/3)$ である。データ要素120が、ある特定の時刻には1つのクラスタに属するけれども、上記3つのクラスタの間で交代する例においては、時刻 t_1 において、 n_1 、 n_2 、及び n_3 は、1、0、及び0の値であってもよい。時刻 t_2 においては、 n_1 、 n_2 、及び n_3 は、0、1、及び0の値であってもよい。 n_1 、 n_2 、及び n_3 の値は変化する場合があるけれども、それらの値の合計は常に1であるべきである。

30

40

【0106】

上記銀行業務の例においては、曖昧に識別された取引先企業へのリスクの代替の割り付けに基づいて、各取引先企業の最大リスク及び最小リスクを知ることにより、何れの任意の取引先企業についても怒りうるリスクのより完全な概念が得られ、不確実な既知知識が明らかとなる。将来又は曖昧さの有望な解消についての現行の確信は、メンバシップの確

50

からしい確度を使用して暮らす他のメンバに重み付けを行うことによって組み込むことができ、これらの重み付けを時間と共に見直して、変化する既知知識を反映させてもよい。

【0107】

ファジーな結合の操作により、2つ以上のデータセットが共通の厳密なキーを共有していない場合でも、それらを組み合わせることができる。例えば、住所が逐語的には同一では無い場合でも、異なるデータベースからの顧客所帯記録を住所について結合することができる。1つのデータセットからの住所が不十分又は不正確である場合、候補整合となる第2のデータセットが複数存在する場合がある。この可能性に、ファジーな結合は適合する。

【0108】

ソート操作は、キーによってレコードを順序付け、レコードのグループに作用する（ロールアップ又は結合等の）キーに基づく操作に先立って使用されることが多い。個々のレコードが複数のクラスタのメンバである（と見込まれる又は実際にそうである）場合に、ファジーなソートを使用して、ファジーなロールアップ等の操作に先立ってレコードを順序付けてもよい。ソート順序の観念及び並べ替えの作用は、複数のクラスタの曖昧なメンバである個々のレコードを複製し、それらを最終的な順序付けにおいてそれらの関連するクラスタの各々に位置付けることによって拡張される。

【0109】

ファジーなデータ操作は、キーを厳密に整合させることに基づくキーグループの代わりにクラスタを使用するという点において、従来のデータ操作とは異なる。クラスタは、キーが *Leslie* である場合に *Lesley* を取り出す上記例を含む。また、クラスタは、*John Smith* はマーケティング部門において半分の時間しか働いていないので彼は厳密には当該部門には属していないにもかかわらず、彼を当該部門に分類する例をも含む。

【0110】

図2Bは、例示的なファジーなデータ操作を図解する。この例において、ファジーなデータ操作180は、キー160に作用し、データセット150を取り出す。キー160は従来のキーである。取り出されたデータセット150は、データ要素151、データ要素152、データ要素153、データ要素154、及びデータ要素155、の5つのデータ要素を含む。これら5つのデータ要素は、キー160に整合しない。しかしながら、それにもかかわらず、それらは上記データ操作によって取り出される。これが、ファジーなデータ操作が従来のデータ操作と異なるところである。キーが与えられると、従来のデータ操作は、そのキーに厳密に整合するデータを取り出す。しかし、ファジーなデータ操作は、そのキーに厳密には整合しないデータを取り出すことができる。

【0111】

クラスタの定義において基本的なこととして、重要なファジーなデータ操作は、異なるレコードにおけるデータの比較である。比較テストを使用して、各クラスタにどのレコードが属するかを特定する。アレンジによっては、比較テストは、各レコードから取得された、選ばれたフィールド値の採点された関数であり、2つのデータ（キーは1つのデータである）の間の定量化された差異を距離としてもよい。

【0112】

(a) 2つのデータ間の距離

2つのデータ間の差異は直感的に単純であることが多い。例えば、*Leslie* という名前と *Lesley* という名前との間の差異は明らかであり、フル・タイムの従業員とパート・タイムの従業員との間の違いは明白である。しかしながら、2つのデータ間の差異を定量化又は測定するのは必ずしも単純ではない。ここで、2つのデータ間の距離を測定するのに使用することができる2つの方法について簡潔に検討する。データ間の差異を定量化する他の方法を下述する原理に基づいて容易に開発することができることが理解されるべきである。ファジーな整合手法及び距離測定の変更例は、例えば、米国特許出願公開第2009/0182728号に記載されている（引用により、本明細書に組み

10

20

30

40

50

込まれる)。

【0113】

(1) 2つの単語間の距離

(例えば、所定の文字集合から形成される) 2つの単語の間の距離を測定する方法(「編集距離」と称されることが多い)には、一方の単語から他方の単語になるのに何回の文字操作がかかるかをカウントすることが含まれる。この例においては、1回の文字操作には単一の文字が含まれる。文字は、様々な方法の何れでコード化されていてもよい。例えば、文字は、シングル-バイト若しくはマルチ-バイトのコード化又は文字集合における文字を表現するのに使用されるコード-ポイントの何れを使用しても、コード化することができる。レーベンシュタイン(Levenshtein)編集距離は、一方の文字を他方の文字に変えるのに必要とされる文字の挿入、削除及び置換の数をカウントする。

10

【0114】

レーベンシュタイン編集距離及びその変形の限界は、それらはオンライン・ファジー整合文脈において(即ち、以前に見たことが無い問合せ単語を有しており、既存の参照集合において整合する変形を見出したい場合に)使用することができないという点である。変形を演算するための削除アルゴリズム(例えば、米国特許出願公開第2009/0182728号。引用により、本明細書に組み込まれる)を代わりに適用することもできる。この方法において、単語間の距離は、各単語から整合する単語に至るまでに必要とされる削除の数をカウントすることによって特定される。図2Cは、LeslieとLesleyとの間の削除距離を如何にして演算するかを示す。操作102により、「Leslie」から「i」が削除されて「Lesle」が得られる。操作104により、「Lesley」から「y」が削除されて「Lesle」が得られる。「Leslie」と「Lesley」との間の距離は1+1(各単語から1回の削除。あるいは、これらの単語の一方のみ、1回の削除と1回の挿入が作用する)。

20

【0115】

アレンジによっては、削除される文字の位置及び相対値を比較することによって、より精巧な採点を行うことができる。これにより、異なる種類の変化に異なる重み付けが適用される、重み付け採点が可能となる。例えば、置換は転位よりも重要ではないとすることもでき、又は「n」による「m」の置換は「k」による「m」の置換よりも重要ではないとすることもできる。

30

【0116】

上記削除アルゴリズムを、以下の方法により、参照用データセットのオンライン・ファジー検索に使用することができる。参照用データセットにおける各単語から1つ以上の文字を(必要と考えられる回数まで)削除することによって得られる全ての単語形成することにより、参照用データセットから削除辞書を構築する。(単語の長さに伴って削除の回数が増えて、より大幅な変形が可能となる場合がある。)元の単語及び削除された文字の位置の両方を、削除の結果として生ずる単語と一緒に記録する。検索がなされる場合、問合せ単語を処理して、1つ以上の文字を削除することによって得られる全ての単語を構築する。これらの単語の各々を上記参照削除辞書において調べて、対応する元の単語を見付け出す。(削除の位置の記録は整合を採点するのに使用される。)次いで、整合した元の単語を、データセットにおける通常の厳密な検索/ルックアップにおいて使用することができる。繰り返すと、問合せ単語が参照用データセットには現れない変形である場合でさえ、この方法は役に立つ。

40

【0117】

もう1つの例は、「Corp.」と「Co.」との間の距離である。「Corp.」から「Co.」にするのに、一方の単語における2回の削除(文字rの削除及び文字pの削除)が必要とされる。故に、2つの単語の間の距離が、各単語について整合する単語を得るのに削除操作が(最少で)何回必要であるかで定義される場合、「Corp.」及び「Co.」は同じ単語「corporation」の2つの互換性のある省略形であるけれども、「Corp.」と「Co.」との間の距離は2+0であるとする事ができる。ユ

50

ーザによって入力されたデータが「Corp .」を使用する一方でデータ操作によって使用されるキーが「Co .」を使用する場合、単語を厳密に整合されることに依存する従来の方法は満足の行く結果を与えないであろう。例えば、ファジーなデータ操作において使用されるキーがABC Co .である場合、キーに厳密に整合するデータ入力のみを取り出す従来の検索データ操作はABC Corp .を与えないであろう。ファジーな検索は、キーの特定の距離（例えば、 $2 + 0$ 又はそれより良好）以内にあるデータ入力を返すように構築することができる。このようなファジー検索においては、ABC Co .というキーに対する整合として、ABC Corp .を返すことができる。

【0118】

あるいは、これら2つの単語は同義語として互換性があるので、「Corp .」と「Co .」との間の距離をゼロ(0)として定義することができる。ユーザ指定の同義語を含むデータ入力を返すようにファジー検索を構築することができる。この例は、ファジーな操作が扱う必要があるかもしれない複雑さを紹介している。

10

【0119】

上記例においては、挿入又は削除等の操作に基づいて、挿入及び削除を両方とも1回の操作としてカウントして、距離を演算する。他のアレンジにおいては、重み付け操作に基づいて距離を演算することができる。重み付けを使用して、別のタイプの操作（例えば、削除）に対して、あるタイプの操作（例えば、挿入）にバイアスをかけることができる。あるいは、重み付けを使用して、別の個々の操作に対して、ある個々の操作にバイアスをかけることができる。例えば、スペースの削除に対応する操作に文字zの削除に対応する操作よりも少なく重み付けして、スペースの脱落は一般的なミススペル間違いであるのに対し、英単語におけるzの挿入は、ことによると、ミススペル間違いというよりも、2つの英単語の間での真の差異であるという事実を反映させてもよい。

20

【0120】

例えば、「sunshine」と「sun shine」との間の距離は、スペースの1回の挿入である。「zinc」と「Inc」との間の距離は、文字zの1回の挿入である。距離の計算において個々の操作が重み付けされない場合、これら2つの距離は等しい(1操作)。1操作の距離以内であれば何れの整合をも返すようにファジーな検索操作が構築されている場合、キー「sunshine」による検索は「sun shine」を返すであろうし、キー「Inc」による検索は「zinc」を返すであろう。

30

【0121】

しかし、重み付け操作を使用する場合、これら2つの距離（「sun shine」と「sunshine」との間の距離、及び「zinc」と「Inc」との間の距離）を異なるようにすることができる。例えば、スペースの挿入を0.5の係数によって重み付けして、スペースの挿入が、タイプミスによって、より起こり易いという事実を反映させてもよい。文字zの挿入を1の係数によって重み付けして、余分な文字zは手違いによって追加され難いという事実を反映させてもよい。図2Dは、「sun shine」と「sunshine」との間の距離が0.5操作である一方、「Zinc」と「Inc」との間の距離が1操作であるように操作を重み付けした場合を示す。

【0122】

0.5文字操作の距離以内であれば何れのキーの整合をも返すようにファジー検索のデータ操作が構成されるアレンジにおいては、キー「sunshine」の検索は「sun shine」を返すであろう。しかし、キー「Inc」の検索は「Zinc」を返さないであろう。

40

【0123】

アレンジによっては、より巧妙な重み付けオプションを定義してもよい。

【0124】

(2) 英国郵便番号(British postal code)間の距離

ファジーな整合が有用なもう1つの用途は、同じ所帯について重複レコードを含む会社の顧客住所データベースで役に立つ。同じ所帯についての複数の入力は、当該所帯に関連

50

する郵便番号における誤植又は当該所帯に関連する名前のミススペルによって起こる場合がある。

【 0 1 2 5 】

起こり得る誤植には、スペースの脱落又は挿入、文字の脱落又は挿入、及び文字のミスタイプを含む。同じ郵便番号においてユーザが2つの誤植をすることはあまり起こらないけれども、珍しくはない。同じ郵便番号においてユーザが3つの誤植をすることはあまり起こらないけれども、不可能ではない。

【 0 1 2 6 】

図3は、郵便番号における誤植によって起こり得る重複レコードを示す。顧客住所データベース300において、John Smithという名前の下に、John Smith ox26qt、John Smith ox26qt、John Smith ox26qy、John Smith ox26qy、John Smith ox26qx、及びJohn Smith ox27qyという、6つの入力が存在する。レコードの全てのペアの間の距離が、レコードのペアを繋ぐ線の隣に表示されている。

10

【 0 1 2 7 】

真正のレコードにおける郵便番号から1+1の削除距離以内の郵便番号を含む何れのレコードも、偽のレコード(誤って入力されたレコード)である可能性が最も高く、上記真正のレコードの重複として扱われると、会社が決定したと仮定する。更に、検索キーの1+1削除距離以内にある全てのレコードを求める検索をファジーな検索として当該会社が定義したと仮定する。

20

【 0 1 2 8 】

ある単語が別の単語の指定の距離以内にある場合、前者は後者の変形である。本例においては、指定の削除距離は1+1(各単語から1回の削除)である。顧客住所データベース300における各郵便番号間の距離情報は図3に列挙されている。図4に示されているように、図4に基づいて、各レコードについての変形を特定することができる。

【 0 1 2 9 】

図4は可視化ツールであり、独特のシェードのボックスで各レコードを表現し、各レコードのボックスをその変形のボックスと重ね合わせることによって形成される。例えば、レコードB、C、及びEはレコードAの変形であるので、レコードAのボックスはレコードB、C、及びEのボックスと重なる。レコードA及びFはレコードEの変形であるので、レコードEのボックスはレコードA及びFのボックスと重なる。

30

【 0 1 3 0 】

場合によっては、どのレコードが真正のレコードであるのかを会社が知っているかもしれないし、他の場合においては、どれが真正であるのか会社は知らないかもしれない。

【 0 1 3 1 】

最初の例において、真正のレコードが「John Smith ox26qy」であるということを会社は知っている。検索キーとして「ox26qy」を使用するファジー検索を実行すると、以下の2つのレコードが取り出されるであろう:「John Smith ox26qt」及び「John Smith ox26qy」。この会社は、真正のレコード「John Smith ox26qy」の重複として、同じクラスタにおけるこれら2つのレコードを扱うであろう。この会社は、これら2つの重複を除去したり、又は当該3つのレコードに共通キーを割り付けることによって、それらをグループ化したりすることを決定することができる。このグループは、ファジーなクラスタの例である。

40

【 0 1 3 2 】

ファジーなクラスタ化は、厳密には整合しないけれども、互いに特定の距離以内にあるキーを有するデータを一緒にグループ化するデータ操作である。ファジーなクラスタ化は、上記例に示されているように、ファジー検索と関係している場合がある。上記の場合におけるように真正のレコードが知られている場合、ファジー検索は、真正のレコードの指定の距離以内にあるデータを取り出す。次いで、取り出されたデータはファジーなクラス

50

タを形成する。

【0133】

どのレコードが真正のレコードであるのかを会社が知らない場合においては、例えば、「John Smith ox2 6qt」及び「John Smith ox2 6qy」の両方が真のレコードであるかもしれず、レコードを如何にして一緒にグループ化するかにあつての先験的助言は存在しないので、単純なファジー検索によつて、どのレコードが互いの重複であるのかを算定して、ファジーなクラスタを作り出すことはできない。このような場合においてファジーなクラスタを生成するのに採用することができる2～3のアプローチにつき、セクション(b)において詳細に説明する。

【0134】

(3) 定量化された差異の他の例

2つのデータの間の距離は、2つのデータの間の定量化された差異の一例である。2つのデータの間の差異を違ふ方法で定量化することもできる。

【0135】

アレンジによつては、整合するペアの間の類似性に基づいて当該ペアを採点するように採点システムを發展させてもよい。次いで、上記ペアの間の定量化された差異を、正規化された整合スコアの補数(complement)として定義することができる。

【0136】

ある事象(例えば、訴訟)の結果が保留されているために不確かさが存在する場合、1つのカテゴリ又はもう1つに属する1つのデータの可能性を使用して、当該データと当該カテゴリを表現するキーとの間の距離を定量化することができる。当該データとカテゴリを表現するキーとの間の距離は、カテゴリが2つしか無い場合に当該データ当該カテゴリに分類されるであろう可能性の補数として、又は、より多くのカテゴリが存在する場合に当該データ当該カテゴリに分類されるであろう可能性の共役(conjugate)として、定義することができる。

【0137】

(b) 変形関係及び変形の(ファジーな)結合

結合操作においてペアとなるデータ要素と比較されるべき、それぞれのデータセットにおけるそれぞれのデータ要素からのオブジェクトは、1つのデータ又は複数のデータの組み合わせとして定義することができる。リレーショナル・データベースにおけるテーブルの行において、オブジェクトは、列における値、値の一部(例えば、サブストリング)、又は2つ以上の列からの値の組み合わせであってもよい。フィールドを含んでなる一連のレコードからなる、単層ファイル・データセットにおいては、オブジェクトは、1つのフィールドにおける値、1つのフィールドの一部又は2つ以上のフィールドの組み合わせであってもよい。文書において、これは、テキストの断片又はテキストのバラバラの断片の組み合わせであってもよい。

【0138】

オブジェクト{k}の集合Sを考察する。Sにおける各オブジェクトkは、変形{v}と呼ばれる、変形オブジェクトの(ことによると、空の)関連する集合を有する。当該関係である下式

【0139】

【数1】

$$k \sim v$$

【0140】

は、「vはkの変形である」と読まれる。アレンジによつては、関数s(k, v)の下での2つのオブジェクトのスコアが閾値T未満である場合、それらは変形であると特定される。

【0141】

10

20

30

40

【数 2】

$$s(k,v) < T .$$

【0 1 4 2】

(採点関数によっては、代わりに、閾値を超えとした方が便利である場合がある。) オブジェクト間の距離は、(上記において考察した文字列についての編集距離又は削除距離の場合、) 単語又は句を比較するための採点関数を構築するための基礎として使用することができる。

【0 1 4 3】

上記変形関係は同値関係(即ち、対称的であり、推移的性質($k \sim k'$, $k' \sim k'' = > k \sim k''$))を有する)である必要は無いけれども、時として同値関係である。変形関係は、たとえ同値関係でなくとも、対称であると仮定される。

10

【0 1 4 4】

【数 3】

$$k \sim v \Rightarrow v \sim k,$$

【0 1 4 5】

即ち、 v が k の変形である場合、 k は v の変形である。

【0 1 4 6】

2つ(又はそれ以上)のデータセット A 及び B の厳密な(内部)結合は、同一のオブジェクト(A における k_A 、 B における k_B)を含む、レコード(行、文書等)のペアリングとして定義される。

20

【0 1 4 7】

【数 4】

$$k_A = k_B.$$

【0 1 4 8】

上記オブジェクト k_A 及び k_B はキーと呼ばれる。

【0 1 4 9】

変形(「ファジーな」)内部結合は、2段階で定義される。まず、レコード(行、文書等)等のデータ要素の暫定的なペアリングが作られる。1つのアレンジにおいて、 A における k_A は、 B における v_{Bn} とペアとなる($k_A \sim v_{Bn}$)。次いで、 k_A 及び v_{Bn} と関連付けられたレコードのペアが評価され($E(k_A, v_{Bn})$)、どのレコードのペアを保持するかが特定される。(厳密な場合においては、全てのペアが保持され、整合ステップ及び採点ステップが単一の比較に融合する($k_A = k_B$))。評価操作は、一般に、ペアリングに使用されるオブジェクト以外の、ペアとなったレコードにおける更なるオブジェクトの比較を含む。アレンジによっては、評価操作はスコアを生成し、このスコアは、整合を識別するには整合閾値を超えなければならない。半結合及び外部結合は、整合するレコードが全く見付からない(又は保持されない)場合に、対抗するレコードに対してヌル値が指定される場合に例えて定義される。

30

40

【0 1 5 0】

最も単純な暫定的ペアリングは下式によって与えられる。

【0 1 5 1】

【数 5】

$$k_A \sim v_{Bn},$$

【0 1 5 2】

即ち、 B における k_A の変形の集合である。このペアリング(「整合」)ステップは、提案されたペアリングが保持されるべきか否かを特定する評価(「採点」)ステップによって補足される。

50

【 0 1 5 3 】

上記変形ペアリング ($k_A \sim v_{Bn}$) への拡張の階層が存在する。第1の一般化は、下式によって与えられる更なるペア (k_A, v_{Bnm}) を付加するによって、上記ペア (k_A, v_{Bn}) を拡張することである。

【 0 1 5 4 】

【 数 6 】

$$k_A \sim v_{An}, v_{An} \sim v_{Bnm}.$$

【 0 1 5 5 】

即ち、 k_A は、Aにおける k_A の変形のBにおける変形とペアとなる。変形関係が同値関係ではない場合、これは、Bにおける要素のより大きい集合に到達する。ここで留意すべきは、この操作は対称的ではない (k_A に到達することができないオブジェクト v_{Bnm} がBに存在する場合がある) ということである。即ち、下式が成立しても、

【 0 1 5 6 】

【 数 7 】

$$v_{Bnm} = k_B, k_B \sim v_{Bi}, v_{Bi} \sim v_{Aij},$$

【 0 1 5 7 】

$v_{Aij} = k_A$ は成立しない。これは、Bにおける k_B の何れの変形も、変形として k_A を有することを必要としないためである。せいぜい、 k_B は、その変形の1つとして k_A の変形を有することを要求される程度である。

【 0 1 5 8 】

変形の変形への更なる拡張等が可能である。特に、下式が成立するペア (k_A, v_{Bnmj}) による (k_A, v_{Bn}) の拡張は、

【 0 1 5 9 】

【 数 8 】

$$k_A \sim v_{An}, v_{An} \sim v_{Bnm}, v_{Bnm} \sim v_{Bnmp},$$

【 0 1 6 0 】

以下の意味において対称的である。(上記操作によって k_A とペアとなる) Bにおける要素 k_B を仮定すると(即ち、 n, m, p によっては、 $k_B = v_{Bnmp}$)、下式が成立する要素 $v_{Aijl} = k_A$ が存在する。

【 0 1 6 1 】

【 数 9 】

$$k_B \sim v_{Bi}, v_{Bi} \sim v_{Aij}, v_{Aij} \sim v_{Aijl}.$$

【 0 1 6 2 】

言い換えれば、同じ変形性号手順を逆に適用すると、逆転したペアが含まれる。Aにおけるオブジェクトから到達されるBにおける全てのオブジェクトは、同じ手順により、今度はAにおける元のオブジェクトに到達することができる。

【 0 1 6 3 】

3つ以上のデータセットへの拡張は、データセットをペア毎に結合し、結果として得られるペアのカルテシアン積をとることによって定義してもよい。A、B、及びCを下式のように結合する。

【 0 1 6 4 】

【 数 10 】

$$k_A \sim v_{Bn},$$

$$k_A \sim v_{Cm},$$

$$\Rightarrow (k_A, v_{Bn}, v_{Cm}).$$

10

20

30

40

50

【 0 1 6 5 】

より高次の拡張は、ペア毎のやり方で上記のように（例えば、変形の変形）定義されるより高次の拡張を使用することによって得られる。状況によっては、随意に、BとCとの間の変形関係が要求される場合がある。例えば、幾つかのn、mについて、下式が成り立つことが要求される場合がある。

【 0 1 6 6 】

【 数 1 1 】

$$V_{Bn} \sim V_{Cm}$$

【 0 1 6 7 】

より高次の変形の使用は、BとCとの間の直接の繋がり（当該関係は、勿論、Aによって既に仲介されている）を確立するのに必要とされる場合がある。

【 0 1 6 8 】

上記において考察したように、変形関係の1つの有用なソースは、編集距離によって関係付けられる単語をペアにすることである。単語間で考察される編集距離が1に限定される場合、これにより、データセット内のペアリングの特定の集合が変形として認められる。例えば、「Smith」は、「Smth」、「Smith2」及び「Smyth」を変形として有する。「Smith20」は「Smith」の変形ではないけれども、編集距離1の変形関係は推移的ではないので、それは「Smith2」の変形である。

【 0 1 6 9 】

上記変形結合は、単一の単語又はフィールド全体を変形キーとして使用することができる場合に使用することができる。例えば、データセットを検索することは、変形キーワードを使用する変形結合として明確化することができる。問合せ句は、キーワードの集合に分解され、その各々は、対象となるデータセットからの単語のインデックスにおける、その変形と整合される。当該インデックスは、対象となるデータセットにおける所定のフィールドにある単語を含むレコードについて、当該単語をレコード識別子（キー）とペアにする。対応するレコード識別子のリストを各キーワードの当該インデックスとの変形整合から得て、これらのリストの積集合を得る（intersect）ことにより、1つ以上のキーワードを共有するレコードを見付け出すことができる。返されたレコードのリストは、整合するキーワードの組み合わせにスコアを割り付けることによって順序付けられる。このスコアは、データセットにおける各キーワードの相対頻度（「逆文献頻度」）（対象となるデータセット・レコードにおけるキーワードの位置と比較した、問合せ句におけるそれらの相対的位置（例えば、順序、隣接性（adjacency））、又は問合せ句からの単語の不在）を考慮に入れてもよい。また、キーワードを関連度の他の尺度と関連付けて、採点をより識別力のあるものとすることもできる。

【 0 1 7 0 】

また、変形結合を、単一の単語についてのルックアップとして使用することもできる。例えば、1つのデータセットにおけるある顧客を、ファースト・ネーム、ラスト・ネーム及び住所によって識別してもよい。この顧客を、第2のデータセットにおいて、ラスト・ネームによって探して、当該整合をファースト・ネーム及び住所によって確認してもよい。当該整合手順は、ソースとなるデータセットにおけるラスト・ネームから、対象となるデータセットにおける様々なラスト・ネームの集合を使用して、整合候補の集合を識別し、取り出すというものである。これらの候補は、ファースト・ネーム及び住所について更に比較され、一致の度合いが整合を識別するのに十分であるかどうか特定される。

【 0 1 7 1 】

例えば、ソースとなるデータセットにおけるレコードが、「Paul, Smith, 20 Walker Street」であり、対象となるデータセットにおける整合する変形の集合が（Smith, Smyth, Smithh）である。対象となるデータセットにおける関連付けられるレコードは以下の通りである。

【 0 1 7 2 】

10

20

30

40

50

【表 1】

- 1,Paul,Smith,20 Walken St
- 2,Robert,Smith,1532 East Grove Ave
- 3,P,Smyth,19 Western Ave
- 4,Pal,Smithh,20 Walker Street

【0173】

確証アルゴリズムは、レコード1及び4が整合となるのに十分に近いことを見出すことができる。これらのレコードを、ルックアップ（検索）によって、又は変形結合において、返すことができる（ここで、2つのデータベースは互いに対して編成されている）。 10

【0174】

あるいは、ことによると、ソースとなるデータセットにおいて、元の「Smith」が、対象において整合するレコード「5, p, Smith 20, Walker Street」を有する「Smith 2」という変形を有する。

【0175】

「Smith 20」は「Smith」の直接の変形ではないけれども、ソースとなるデータセットにおける変形「Smith 2」からは到達することができる。

【0176】

変形結合のもう1つの使用は、クラスタ化に先立つスーパークラスタの定義である。これについては、クラスタ化が定義されてから、後に考察する。 20

【0177】

(c) クラスタ化及び部分的メンバシップ

多くの厳密なキーに基づく操作は、共通キー値を共有する集合にレコードをグループ化することを必要とする。これらの集合は、時として、「キーグループ」と呼ばれる。例えば、ロールアップ操作は、あるキーグループにおけるレコードに亘ってデータを組み合わせたり、統合したりする。カウント、総計、最大値又は最小値、値のベクトル、一意的な値への重複等は全て、ロールアップ操作によって演算することができる。レコードのグループを単一のレコードに集約する操作は何れも、ロールアップ操作として解釈することができる。 30

【0178】

（独立した処理のためのデータ・パーティションにデータが分離される）データ並列処理は、同じキーグループに属する全てのレコードが同じデータ・パーティションに存在することを保証するために、キーに基づく分割に依存することが多い。ロールアップ及び結合等の操作は、順次（非並列）処理における結果と同じ結果を生成するために、このことに依存する。

【0179】

キーグループの集合は、全てのレコードの集合を互いに素な集合へと分割することを構成する。全てのレコード（オブジェクト）は唯一のキーグループのメンバである。クラスタは、キーグループの観念を重なり合う集合を含むパーティションへと一般化する（メンバシップはキーの厳密な一致によって特定されない）。 40

【0180】

各々が重み付け $w(k, C)$ を有する、オブジェクト k の重なり合うかもしれない集合 $\{C\}$ （クラスタと呼ばれる）の集まりに集合 S を分割（分解）することを考察する。あるオブジェクト k は2つ以上のクラスタ C のメンバであるかもしれない、もしそうである場合、そのクラスタ・メンバシップは、曖昧又は部分的であると言われる。 C における k に割り付けられる重み付け $w(k, C)$ は、 C における k の「部分的メンバシップ」を定量化し、時として曖昧さの尺度と呼ばれる。クラスタ C は、ペア $C = \{(k, w(k, C))\}$ の集合として表示することができる。 $w(k, C) = 0$ である場合、 k は「 C のメンバではない」と言われる。 $w(k, C) = 1$ である場合、 k は「確かに C のメンバである 50

」と言われる。固定された k について、 C を亘っての重み付けの合計は 1 に等しい (S における確かなメンバシップに対応する)。

【 0 1 8 1 】

【 数 1 2 】

$$\sum_C w(k, C) = 1.$$

【 0 1 8 2 】

重み付けの割り付けは、ルール R と関連付けられ、 R によって表示される。所定の集合 S は、概して、クラスタの集まりへの 2 回以上の分割、及び異なるルールの下での各オブジェクトへの 2 回以上の重み付けを認める。一般に、異なるルールに関連付けられる重み付けを組み合わせることはできない。

10

【 0 1 8 3 】

クラスタ $C = \{ (k, w(k, C)) \}$ の補数は、集合 $\{ (k, 1 - w(k, C)) \}$ であると定義される。特に、当該補数は、重み付け 1 を用いて、 C にはないオブジェクトを表す。クラスタ $\{ C \}$ の集まりが S にかからないか、又は C を亘っての k についての重み付けの合計が 1 に等しくない場合、 S における $\{ C \}$ の和集合の補数は、当該集まりに付加されたものと仮定される。

【 0 1 8 4 】

2 つのクラスタを組み合わせる単一のクラスタとし、それにより、それらの重み付けを合計することによって、分割を粗くすることができる。

20

【 0 1 8 5 】

【 数 1 3 】

$$C_1 + C_2 = \{ (k, w(k, C_1) + w(k, C_2)) \}.$$

【 0 1 8 6 】

このプロセスを逆転させて、新しいクラスタの中の各オブジェクトに重み付けを割り当てることによって、クラスタを更なるクラスタに分解することができ、それ故に、新しい重み付けの合計は、元の重み付けに等しい。例えば選択基準を適用した後に、オブジェクトの重み付けを減じることによって、クラスタからオブジェクトを取り去ることができる。

30

【 0 1 8 7 】

状況によっては、重み付けが $0 \leq w(k, C) \leq 1$ を満足する場合、これらの重み付けは「 k が S におけるクラスタ C のメンバである可能性」との解釈を認めることができるけれども、一般には、クラスタの当該定義は非統計的である。負の重み付け及び 1 よりも大きい重み付けの可能性は排除されないけれども、固定された k についての重み付けの合計は 1 でなければならない。

【 0 1 8 8 】

$C_1 + C_2$ 等の組み合わせは、原則として、形成することができるけれども、それらは、あるオブジェクトについてあり得る最大の重み付けの値によって反映されるように (例えば、この場合は 2)、オブジェクトの複数のコピーを含むクラスタを構築することに該当する。各オブジェクトについての重み付けの合計は 1 であるという条件は、集合 S において各オブジェクトのコピーが 1 つだけ存在することを想定している。それが当てはまらない場合、重み付けの合計の値は、それ相応に変化する場合がある。一般に、オブジェクトによって重み付けの総計が変動することを排除するものは無い。

40

【 0 1 8 9 】

部分的メンバシップを有するクラスタは、集合における各オブジェクトに重み付けを割り付けるメンバシップ関数を有するオブジェクトの集合として記述することができるので、ファジーな集合の概念に類似している。しかしながら、ここでは、孤立しているファジーな集合には重きが置かれなければならないけれども、分割の要素としてのクラスタに重きが置かれる

50

。特に、重み付けは、分割の性質であって、孤立しているクラスタ内のオブジェクトのものではない。クラスタ内のオブジェクトに割り付けられる重み付けは、他のクラスタに対して起こり得る割り付けによって影響を受ける。焦点は、クラスタ内のメンバシップ関数からクラスタを亘ってのオブジェクトへのメンバシップの割り当てにシフトする。

【0190】

部分的メンバシップを有するクラスタは多くの状況において自然に生ずる。最も単純な言葉によれば、部分的メンバシップは、オブジェクトのクラスタへの割り付けにおける曖昧さの結果である。厳密なキーが存在する場合、オブジェクトが何れのキーグループに属するのかという疑問は存在しない。厳密に一致することを必要としない1つのデータ又はデータの組み合わせにメンバシップが基づく場合、メンバシップ決定はあまり明確ではない。

10

【0191】

以下は、部分的メンバシップに繋がり得る幅広いデータ品質問題の例である。

【0192】

データは、クラスタ化のルールと比べて、本質的に曖昧である。時として、クラスタは、それらの定義が排他的メンバシップを推定していないという単純な理由により、重なり合う。ある会社において2つの部門で働く従業員について考察する。従業員のリストが部門によってクラスタ化されている場合、当該従業員は2つのクラスタに現れるのが正しい（これが事柄の真の状態を反映している）。この場合、部分的メンバシップを設定して、当該従業員が各部門のために働く時間の分配を反映させることができる。このことは次に、別個の部門に同じ識別情報を有する2人の別個の従業員が存在するという誤った結論を導き出す機会を減らす。

20

【0193】

データは不完全である場合がある。フィールドにおける変形の単語が識別を曖昧にする場合がある。住所を所帯に割り付ける際に、1つの住所レコード上の番地が12であり、その通りには12という番地は無い場合がある。代わりに、1、2及び21という番地が存在する場合がある。最初の2つの番地への0.3の重み付け及び最後の番地への0.4の重み付けにより、挿入エラーよりも僅かに大きい転移エラーの可能性を反映させることができる。

【0194】

データは不十分である場合がある。クラスタへのえっていきな割り付けを行うのに必要な情報が欠落している場合がある。例えば、住所レコードを所帯に割り付ける問題を考察する。各々の一意的な番地、通り、市、郵便番号の組み合わせには、一意的な所帯番号が割り付けられる。クラスタ化アルゴリズムは、通りの名前及び市の名前の変形のスペルに寛容である場合があり、それ故に、正しい所帯に割り付けられる全ての住所が同一である必要は無い。しかしながら、住所から番地が欠落している場合、信頼の置ける割り付けを行うのに十分な情報が無い。できるだけ多くの情報を保持するために、不十分なレコードを、入手可能な情報と一致する各所帯に部分的に割り当ててもよい。可能性がある番地が5つ存在する場合、各番地クラスタにおける重み付けは、各番地の等しい確度を反映させて、0.2とすることができる。異なる見地で、会社名によって表示された会計における負債について考察する。銀行は、このデータを統合して、国毎に各会社と関連付けられる未払いの負債を特定することが必要である。これらの会社名の中には、「ACME SERVICES LIMITED (AUSTRALIA)」、「ACME SERVICES LIMITED (CANADA)」及び「ACME SERVICES LIMITED」がある。最初の2つの各々は別個のクラスタに行くけれども、3つ目は、最初の2つの各々に対して同等の整合であり、国の識別子が足りない。最初の2つのクラスタの各々における3つめの会社に0.5の重み付けを行うことにより、会社情報の不十分さを反映させる。

30

40

【0195】

データ又は分類は、本質的に不確かである場合がある。クラスタのメンバシップは、将

50

来の事象の結果に基づいている場合がある。資産及びそれらの価値のリストを含むデータセットについて考察する。資産は、所有者によってクラスタ化されるべきである。しかしながら、特定の資産の所有権についての訴訟が係争中である。その資産を可能性がある所有者の何れかに結び付けることは負け博打 (l o s t b e t) である場合があるけれども、その資産を単純に無視することもできない。

訴訟の予想結果の知見の現状を反映する部分的メンバシップを用いて当該資産を各所有者に割り付けることは、現行の知見と一致した最も公平で有益な当該資産の処分を生ずる。

【 0 1 9 6 】

(d) クラスタ化

クラスタ化は、クラスタ・メンバシップ基準に基づいて、レコードをクラスタにグループ化する行為である。厳密な場合、各レコードからのオブジェクト (キー) が、他のレコードにおける対応するオブジェクトと厳密に整合され、クラスタ (又は「キーグループ」) は共通キーを共有するレコードの集合である。ファジーな場合、クラスタ・メンバシップを、各レコードにおけるオブジェクト間の変形関係によって特定してもよい。(より一般的には、クラスタ・メンバシップ基準があり得る。) データセットにおける全てのレコードを互いに比較しなければならないことを回避するために、スーパークラスタ・キーを使用して、全体集合を部分集合に分割して、1つのスーパークラスタ内のレコードに対して相互比較を制限する。

【 0 1 9 7 】

多くの場合、スーパークラスタは、厳密なキー (例えば、郵便番号) によって定義される。変形結合は、変形オブジェクトを使用してスーパークラスタを定義することを可能とする。例えば、スーパークラスタは、所定の郵便番号の変形である全ての郵便番号を含むレコードの集合として定義することができる。例えば、英国の郵便番号 O X 2 6 Q Y を考えると、変形郵便番号 O X 2 6 Q Y 及び O X 2 6 Q T は両方とも編集距離が 1 の変形であり、後者はそれ自体が有効な郵便番号である。各変形郵便番号からのレコードを有効な整合として認めることにより、郵便番号におけるエラーに寛容なクラスタを生ずることが可能となる。

【 0 1 9 8 】

もう1つのアレンジにおいて、各レコードにおける選ばれたフィールドから (例えば、相対頻度に基づいて、最も長い又は最も重要な単語から) 単語の断片を取り、この断片の変形を使用してスーパークラスタを識別することによってスーパークラスタを形成してもよい。これは、同じクラスタのメンバであるべき2つのレコードについて、それらが特定の単語を共有する可能性が非常に高いけれども、それらの単語が直接的な変形である必要は無く、ましてや等しい必要は無い場合に適している。変形断片をスーパークラスタ・キーとみなすことによって、その単語の残りの部分が変形として許容されるよりも異なっているレコードが認められる。クラスタ・メンバシップを特定するには、各レコードにおけるフル・ワード及び他のオブジェクトの徹底的な比較が必要とされる。

【 0 1 9 9 】

例えば、ドイツ語の通りの名前を比較する場合、 G r a f v o n S t a u f f e n b e r g という通りは、 S t a u f f e n b e r g を何等かの形で吹くことが必要とされる場合がある。サンプル・データにおいて、 S t r a s s e が省略され、先行する単語に連結されて、 G r a f v . S t a u f f e n b e r g s t r 等の入力を生じていることが観察された。各通りの名前における最も長い単語の最初の5つの文字の編集距離が2の変形によって定義されるスーパークラスタは、「 s t a u f 」及び「 s t a u f f 」の両方を含むであろう。 S t a u f f e n b e r g 及び S t a u f f e n b e r g s t r の両方を含むレコードが比較のためにスーパークラスタに含まれ、好適な採点関数がそれらに同じクラスタを割り付けるであろう。一方、最も長い単語の編集距離が2の変形に基づくスーパークラスタは、これら2つの通りの名前を別個のスーパークラスタに区分けするであろう。そして、それらが一緒にクラスタ化されることはできない。

【 0 2 0 0 】

10

20

30

40

50

スーパークラスタを思慮深く選ぶことは、クラスタ化方法の性能及び精度のために重要である。例えば、スーパークラスタが大き過ぎる場合、性能を劣化させ得る、多くの報いの得られない比較がなされる場合がある。あるいは、スーパークラスタが狭過ぎる場合、許容可能な整合が欠落し、精度が劣化する場合がある。

【0201】

(e) 部分的メンバシップ

データ入力オペレータが、新しい顧客をデータベースに追加するために、アプリケーションにおけるフォームに記入していると仮定する。名前をフォームに入力する際に、アプリケーションは、名前の参照リストと対照して、入力を確認する。上述のような削除アルゴリズムによるファジー検索を使用して、アプリケーションは名前の変形スペルを検出し、上記参照リストから他の可能性のリストを返すことができる。オペレータがファースト・ネームのフィールドに「J a m e」を入力すると仮定する。アプリケーションは、アルファベット順の他の可能性の以下のリストを（その名前を含むデータベースにおけるレコードのカウントと一緒に）返すことができる。

【0202】

【表2】

Jaime 250

James 13359

Jamie 339

Jane 9975

【0203】

これらは全て、1回の挿入及び/又は1回の削除だけ、J a m eとは異なっており、候補選択肢である。

【0204】

上記リストのオペレータに対する実用性を向上させるために、曖昧さの尺度を特定するための様々なモデルの何れかを使用して、選択肢を優先順位付けすることができる。曖昧さを定量化するための3つの例示的な方法は、1)均等分配、2)統計的頻度、及び3)誤差モデルである。

【0205】

均等分配アプローチにおいては、各選択肢は等しく起こり得るものとして扱われる。ここで、J a m eが上記選択肢の何れかである確度は四分の一である。選択肢のアルファベット順のリストは、概して、黙示的な均等分配アプローチを示唆する。

【0206】

統計的頻度アプローチにおいては、(データベースのテーブル自体等の)参照集合を各名前についての観測度数として使用する。上記リストを示されているカウントの降順でソートする場合、最も起こり得る修正はJ a m e sであり、続いてJ a n e等である。

【0207】

誤差モデルの第3の方法は、とりわけ、言語及び入力のモードによって、特定の種類の誤差が他の種類の誤差よりも起こり易いという観測結果に基づく。熟練したオペレータによるキーボード入力については、文字を飛ばしたり又は余分な文字を挿入したりするよりも、置換エラーの方がよく起こる場合があるかもしれない。同様に、電話越しに顧客によって与えられるスペルを記録しているオペレータについては、音声学的に類似の文字の名前を含む転記エラーの方が、他の種類のエラーよりも起こり易い可能性がある。何れの場合においても、ここでは、J a n eが最も可能性のある修正であろう。この方法を使用するために、起こり得るエラー及びそれらの相対的重要度を分類するモデルを開発し、適用することができる。このようなモデルは、米国特許出願公開第2009/0182728号において紹介されている、WFS(単語頻度優位性: word frequency significance)ファイルの統計的分析から作り出すことができる。

10

20

30

40

50

【0208】

アプリケーションが顧客住所（ここでは問合せ住所と称する）を含むレコードをマスタ顧客住所テーブルに整合させて、整合が見付かった場合には、存在する所帯キーを取り出し、又はさもなければ新しいキーを作り出すと仮定する。問合せ住所は、マスタ顧客住所テーブルにおける住所と厳密には整合しない場合があるので、ファジーな整合を使用することができる。その上、問合せ住所は不十分又は不正確である場合がある。このことは、2つ以上の実在する住所が問合せ住所に対する整合となる場合があることを意味する。整合の品質を定量化するには、整合の曖昧さの尺度を有することが有用である。

【0209】

例えば、住所は番地を有しておらず、マスタ顧客住所ファイルは同じ通りの住所を有する複数の入力を有している場合がある（番地は無視する）。問合せ住所は特定の町及び郵便番号を有する Lower Street であると仮定する。郵便番号についてのファジー検索は、同じ郵便番号又は変形郵便番号を有する、見込みのある住所レコードのリストを返す。問合せ住所上の郵便番号、町、通り、及び番地のフィールドは、ファジーな整合プロセスの一部として、見込みのある住所の対応する各フィールドと比較され、採点される。この例において、問合せ住所の、通り、町及び郵便番号に厳密に整合する2つのマスタ・レコード（2 Lower Street 及び 3 Lower Street）が存在すると仮定する。各有力候補は問合せレコードとの整合の品質が等しく、問合せレコード上の番地が欠落しているため、既存のデータではこれを改良することはできない。均等分配の尺度の下では、何れの所帯に対しても、実際の整合の確度が等しい。

【0210】

あるいは、番地は埋まっているけれども無効であり、実在する住所の何れにも符合しないと仮定する。問合せ住所は 12 Lower Street であるけれども、参照用のポータル・アドレス・ファイル（ポータル・サービス（Postal Service）から入手可能な全ての有効な住所のリスト）に照らした確認により、その住所を有する家が無いことが示されると仮定する。その郵便番号において整合する住所は、上記のように、2 Lower Street 及び 3 Lower Street である。住所入力についての誤差モデルは、3 に対する 12 の整合よりも、2 に対する 12 の整合の方を選ぶ。これにより、整合の確度の重み付けにバイアスをかけて、2 Lower Street の住所との整合に有利に働くであろう。

【0211】

最後に、問合せレコード上の番地が埋まっていて且つ有効な郵便の宛先である場合、住所入力についての誤差モデルは、当該住所が実在する住所のエラーに対して新しい確度を定量化することができる。

【0212】

(f) データ品質の定量化

曖昧さの尺度は、データ品質を測定する、より広い状況においても適用可能である。企業及び組織は、それらのデータ、特にそれらのマスタ・データの品質について懸念するけれども、現在のところ、最も明らかなデータ品質問題以外のものを定量化することは困難である。上記に示したデータ品質問題の短いリストの中で、幾つかのデータ品質測定システム（例えば、米国特許出願公開第 2005/0114369 号を参照。引用により、本明細書に組み込まれる）は、主に、1つ（データ妥当性）に直接的に取り組む。データは網羅的に目録に載せられ、そのデータ型及びユーザが定義する妥当性の様々な尺度（有効値のリストを含む）に照らして、妥当性がチェックされる。

【0213】

レコードのフィールド内の不十分な入力の証拠は、埋まっていない（空又はヌル）入力の数から推察することができるけれども、これでは、欠落した情報の重要性は定量化されない。マスタ顧客住所リストの場合を考察する。顧客住所入力から市が欠落しているけれども、有効な郵便番号及び通りの名前がある場合、これは何等かの曖昧さを加えるであろうか？又は、ことによると、ポータル・アドレス・ファイルのような参照集合を使用して

10

20

30

40

50

、手元の情報から住所を効率的に完成することができるか？もし住所から番地が欠落していたらどうなるか？残りの住所を何軒の家が共有しているか？幾つかのデータ品質問題はデータ・クレンジングによって解決することができ（欠落している市）、他は解決することができない（欠落している番地）。データに存在する本質的な曖昧さの尺度が必要とされている。

【0214】

各アドレス入力を参照用データセットと比較することにより、当該入力における曖昧さの尺度を演算することができる。曖昧さのレポートには、曖昧さが全く無い入力の割合が含まれる場合がある。曖昧さを有する入力については、レポートは、K個の選択肢を有する入力の数のヒストグラム・プロット（又は、特定の瓶/範囲の選択肢、「変形」とも呼ばれる）を示す場合がある。また、最も大きい曖昧さを有する最初のN個の入力のリストもあるかもしれない（ここで、Nはユーザによって指定される入力数である）。データセット全体について不十分なデータと関連付けられる曖昧さを定量化する1つの要約統計量は、入力当たりの選択肢の数の平均偏差及び標準偏差から構築される。

10

【0215】

曖昧さの統計的頻度を適用して住所の完成の可能性を定量化する場合、興味深い尺度は、a) K個の選択肢を有する入力の数のヒストグラム・プロット、b) 選択肢の頻度の分布のヒストグラム・プロットを伴う、選択肢の最も大きい範囲内のN個の入力のリスト、c) 単一の選択肢への最強の関連を伴うN個の入力のリスト、d) 選択肢の数の平均偏差及び標準偏差である。

20

【0216】

類似の尺度が、曖昧さの誤差モデル尺度に当てはまる。

【0217】

変形入力のデータ品質問題（フィールドにおける入力が意図したものではない）は、無効値の問題及び不十分な情報の問題の両方に類似している。レベルによっては、ある入力が意図したものではないと言うことは、特定の（必ずしも枚角ではないけれども）検証基準に対して無効であるということと断言することである。通りの名前の妥当性は、当該通りの名前を参照用のポータル・アドレス・ファイルに含まれているものと比較することによって定義することができる。あるいは、参照用データセットが無い場合、関連付けられた変形に整合する入力の相対頻度から妥当性を推察することができる。ある入力が相対的に高い頻度で起こり、高い頻度の選択肢が無い場合、それを妥当として解釈することができる。それが相対的に低い頻度で起こり、単一の高い頻度の選択肢が有る場合、それを無効として解釈することができる。複数の高い頻度の選択肢が有る場合、妥当なデータへの修正は、曖昧である場合があり、当該データが欠落していたらどうなっていたであろうかということと同様にして定量化することができる。

30

【0218】

フィールドが有効値の列挙集合の範囲外の値を有すると仮定する。例えば、性別フィールドが、M又はFの代わりにGである。曖昧さの均等分配尺度によれば、当該入力には2つの選択肢があったと言われる。頻度尺度もやはり2つの選択肢を示すけれども、M又はFのいずれかへのバイアスを含む場合がある。選択肢の数に変動が無い、単純な列挙された場合においては、当該データセットについての曖昧さの尺度は、入力当たりの曖昧さの尺度掛ける無効値の割合の積から形成される。

40

【0219】

もし通りの名前にミススペルがあったらそうであろうように、選択肢の数に変動が有る場合、曖昧さの尺度が、当該データに存在する曖昧さを定量化する助けとなるであろう。頻度及び誤差モデル尺度は、閉経のスペルについて最も信頼性のある結果を与えるであろう。不十分な情報と同様に、曖昧さの尺度は、究極的には、当該データセットをクレンジングによってどれくらい向上させることができるか、及び不確かさが未だどれくらい残るかを反映する。

【0220】

50

間違っただフィールドに入力されたデータもまた、同様の方法で定量化することができる。ここで、データ配置が実際に間違っているか否かに関する更なる曖昧さが存在する場合がある。一般的なファースト・ネームの値を有する名字の場合、ある順序付けが別の順序付けよりも正しいか否かがハッキリしない場合がある。一般的なファースト・ネームが名字フィールドにどれくらい頻繁に存在するか（又はその逆）を知ることは、間違っただ順序付けされた名前の問題が如何に深刻であるかを明確化するのに役立つ。正確な参照（又は他の検証手段）が入手可能である場合、エラー率の尺度を得ることができる。100,000個の入力のデータセットのうち、一般的なファースト・ネームである名字が500個あるけれども、これらのうち25個だけが実際に間違っただ順序付けされている。そのときのエラー率は、 $25 / 500 = 1 / 20$ である。測定されたエラー率が無い場合でも、上記問題に対して脆弱な入力がある場合、 $500 / 100,000 = 1 / 200$ しかないことを知ることは、そのデータにおける信頼度を向上させる。

10

【0221】

名前フィールドについての他の興味深いチェックは、a)名及び姓（名字）が一般的なファースト・ネームである場合、b)名が一般的なラスト・ネームであり、姓（名字）が一般的なファースト・ネームである場合、c)名が一般的なラスト・ネームであり、ラスト・ネームも同様である場合、において知られるべきであろう。名前の参照用データセットから導かれる曖昧さの頻度尺度（ことによると、データセットそのもの）を使用して、順序付けが正しい確度を演算することができる。例えば、David Paulという名前について考察する。Davidがファースト・ネーム又はラスト・ネームである可能性（及びPaulについても同様に）から、Davidがファースト・ネームである確度を演算することができる。

20

【0222】

例えば、住所におけるフィールド間での並べ替えは、選ばれた参照用データセットと比べて、間違いである場合があるけれども、順序付けを指定する標準が弱い又は存在しないの何れかであるため、正確に言えば、間違いではない。ここで、参照用データセットと比較した特定の住所の順序付けの曖昧さの尺度を使用して、異なる順序付けは、データへの関連付けの曖昧さを導入しないので、深刻なデータ品質問題ではないことを示すことができる。このことは、データ品質を向上させるために、どこに労力を注ぎ込むべきかの決定に重要な情報である。

30

【0223】

フリー・テキスト・フィールドに存在する参照用データの量を定量化するために、当該フリー・テキストを、選ばれた参照用データセットに照らしたファジー検索において使用される単語に分解することができる。例えば、製品番号が請求書のコメント・フィールドに格納されていると会社が予想していると仮定する。当該コメント・フィールドにおける各単語を使用して、製品番号テーブルに照らしてファジー検索を行うことにより、当該コメント・フィールドに存在する製品番号の数を見出すことができる。より一般的には、各単語をWFSファイル（他の特許を参照）に照らして検索して、その単語又はその変形が現れるフィールドの全てを、どれほどの配分率であるかと共に、特定することができる。これにより、あるフィールドにおけるデータの他のフィールドにおいて見られるデータに対するファジーな相互相関が与えられる。（このアプローチを使用して、間違っただフィールドに配置されたデータを識別することもできる。）

40

【0224】

多くの参照整合性問題は全て、定量化することができる。まず、リンクされたデータが実は正しくリンクされていない厳密なキーのペアの割合として、確証尺度を定義することができる。この種の確証は、比較可能なフィールドがリンクの両側に存在する（又は、更なる結合を通してリンケージの点まで繋ぐことができる）場合にのみ行うことができる。概して、これは、データベース間に結合キーが確立されている2つ以上のデータベースに保持された（顧客の名前及び住所等の）類似のフィールド間での比較である。この状況は、データ・ウェアハウスに保持されている比較データが、データ・ウェアハウスを占める

50

ソース・システム中のデータと比較される場合に起こる。異なるソース・システムが相反する情報を保持している場合、又はそれらが一貫性無く更新されている場合、データ・ウェアハウスも同様に1つ以上のソース・システムと相容れない場合がある。データ・ウェアハウスをそのソースとの一貫性について妥当性を確認することにより、企業データ品質についての新しく重要なチェックが提供される。

【0225】

第2のチェックは、リンクの向こう側にレコードが無いリンクである、欠落リンク及び孤児リンクを探すことである。ファジー検索により、リンクされるべき別のデータが有る(欠落リンク)のか又は無い(孤児リンク)のかを特定することができる。これらの条件の各々の割合の尺度が重要である。リンクの向こう側のレコードへの整合が曖昧である場合、曖昧さの尺度は、それを定量化することができる。このことは、このことを一意的に行うことができるリンクを再投入し、どのリンクが曖昧であるが故に更なる調査を必要としているのかを識別する、データ・クレンジング操作の基礎を形成する。(以下、部分的メンバシップの可能性について考察する。)

10

【0226】

2つのデータベースが関係する情報を含むけれども、それらの間にキー関係が存在しない場合、データセット間でのファジーな検索又は結合により、これらのデータセット間に予想されるリンクが見付け出されるであろう。予想される各リンクの曖昧さの尺度は、これらのデータセットの間でのマッピングが如何にクリーンであるかを示すであろう。このことは、例えば、2つの会社が合併する際に、それらのマスタ参照用データ(例えば、顧客の名前及び住所)を組み合わせる場合に、非常に有用である。同様に、それを使用して、企業の異なる部分の参照用データをマージすることができる。このことは、マスタ・データ管理ソリューションの準備の初期段階で重要であろう。マスタ・データ管理ソリューションの創作におけるギャップ分析の部分は、異なるシステムによって使用される現存する参照用データセット間での調整の品質を決定する。システム間での初期調整は、この分析の副産物である。次いで、曖昧さの尺度を有することが、これらのシステムを調和させるのにかかる更なる労力を定量化するのに役立つ。

20

【0227】

(g) クラスタ及び部分的メンバシップ

(1) ファジーなクラスタの創作

前に述べたように、如何にして要素を纏めてグループ化するかについての先験的助言が存在しない場合、原理又はアルゴリズムを使用してグループを識別する。実際の場合、どの要素が他の要素を引きつけてグループを形成する中核的要素として役に立つかは明らかではないことが多いので、このことは実用的見地から重要である。例えば、顧客住所データベースにおける重複レコードの上記例において、どのレコードが真正のレコードであるかを会社が見分けるのは時として不可能である。以下の考察により、それらの中で幾つかのデータが互いの仲間であるとみなされる、レコードのファジーなクラスタを形成するのに使用することができる幾つかのアルゴリズムが提案される。

30

【0228】

上述の顧客住所の例において、「John Smith」と関連付けられる6つのレコード：(A) John Smith ox2 6qt; (B) John Smith ox2 6qx; (C) John Smith ox2 6qy; (D) John Smith ox2 7qy; (E) John Smith ox2 6qt; (F) John Smith ox2 6qy、が存在する。どのレコードが実在の所帯に該当するのかを知らず、上記レコードを、各クラスタが実在の所帯を表現する2つ又は3つのクラスタにグループ化することに会社が感心を抱く場合がある。この方法において、会社は、偽の郵送先住所に送られる郵便物を減らすことにより、郵便物の量を減らすことができる場合がある。

40

【0229】

クラスタを作り出すのに使用することができる1つのアルゴリズムは、指定された距離

50

内にある最大数の要素を含む最大の互いに素な部分集合を見付けることである。このアプローチについて、図5を参照しながら説明する。含まれるステップは、フロー・チャートに図解されている。図5もまた、上記顧客住所データベースの例を使用して、当該アルゴリズムを確認する。各ステップからの結果は、対応するステップの右側に示されている。この例において、指定された距離は2操作である。

【0230】

図5におけるフロー・チャートに言及して、最大の互いに素な部分集合からクラスタを作り出す最初のステップは、各要素について、その要素の変形をカウントすることである(502)。上記において定義したように、要素の変形は、その特定の要素から指定された距離内にある要素である。顧客住所の例においては、レコードAについては、2操作の距離内にあるレコードが3つ(B、C、及びE)存在する。レコードBについては、2操作の距離内にあるレコードが2つ(A、C)存在する。Cについては4つ、Dについては1つ、Eについては2つ、そしてFについては2つ存在する。

10

【0231】

次いで、最大数の変形を有する要素を選択し(504)、当該要素及びその変形を、当該要素によって表示されたグループにする。顧客住所データベースの例においては、レコードCが最大数(4)の変形を有する。レコードC及びその変形(A、B、D、F)が第1のクラスタを形成する。

【0232】

次に、全ての要素の集合から、最大のグループの要素を除去する(506)。顧客住所の例においては、これにより、レコードEだけが残る。

20

【0233】

次いで、残っている要素において、最大数の変形を有する要素を見付ける(508)。このステップにより、第2のクラスタが生ずる。顧客住所の例においては、第2のクラスタは、その中に1つの要素Eしか有していない。

【0234】

全ての要素がクラスタにグループ化されるまで継続する(510)。顧客住所の例においては、全ての要素がそのグループを見付けてしまったので、更に続ける必要は無い。

【0235】

このアルゴリズムにより、顧客住所の例において、2つのクラスタ(A、B、C、D、Fからなるグループ、及びEのみからなるグループ)が生成する。当該会社は、各グループに含まれるレコードを互いの重複として扱うことができ、これらのレコードを統合して、郵便物の量を減らすことができる。

30

【0236】

上記アルゴリズムに若干の調節を加えてもよい。例えば、レコードA及びFは両方とも、C及びEから同じ距離にある。上記アルゴリズムにおいては、レコードA及びFをCの重複として割り付けることは人為的であり、レコードA及びFがEよりもCに近いことを必ずしも示さない。

【0237】

1つの調節は、クラスタの表現についての不確かさを注記することである。例えば、C5-2という表現を使用して、レコードC及びその変形を含むCのクラスタを、レコードの総数を示す5及び不確かさを示す-2で表現することができる。E1+2という表現を使用して、レコードE及びその変形を含むEのクラスタを、そのグループにおける1及び不確かさを示す2で表現することができる。クラスタの正の不確かさは、他のところにグループ化されている要素が、このクラスタに属しているかもしれないことを反映している。クラスタの負の不確かさは、このクラスタにおける要素が別のグループに属しているかもしれないことを反映している。

40

【0238】

もう1つの調節は、A及びFをEのクラスタに加えることであってもよい。従って、グループCはレコードA、B、D、Fを有し、グループEはレコードA及びFを有する。し

50

かしながら、レコード A 及び F は 2 つのグループに属するため、全てのグループにおけるレコードの総数は 8 であり、レコードの真の総カウントよりも 2 つ多い。真のカウントを維持するために、部分的メンバシップを使用してもよい。

【 0 2 3 9 】

ファジーなクラスタを構築する第 2 の方法は、会社名等の、複数の単語のフィールド上のデータをクラスタ化する際にふさわしく、この場合、単一の単語ではなく句（又は、英国の有文番号のようにフィールド全体）を採点することによってレコード間の変動が評価される。句の採点は、単語の変形スペルのみならず、単語の順序、欠落した単語、及び詳細な単語の連なりを変化させる単語間への挿入をも考慮に入れる。例えば、Bank of America という会社名について、識別され、区別される必要がある名前の変動の 4 つのタイプを如何に説明する。

【 0 2 4 0 】

【表 3】

- 1) Bank of Amrica (fuzzy match of a word)
- 2) America Bank (word order, missing word)
- 3) Bank America (missing word)
- 4) Bank of South America (inserted word)

【 0 2 4 1 】

句に基づくクラスタ化の例として、マスタ顧客リストにおける同じ法人に属する会計の全てを識別することを銀行が試みると仮定する。法人は、会社名、住所、及び（存在する場合は）会社登録番号によって識別されるであろう。会社名は法人と非常に相関性があり、常に入力されているので、クラスタ化に使用される主フィールドは会社名である。住所は、類似の名前を偶然有する会社を区別するのに使用される第 2 のフィールドである。会社登録番号は法人の識別において最も確実であることが期待されるけれども、単独で使用するのに十分なほど入力されてはいない。

【 0 2 4 2 】

ファジーなクラスタ化操作は、元のデータセットを、クラスタ・メンバシップについての全ての要素の比較を可能とするのに適切な大きさの、より小さい部分集合に分割するスーパークラスタ・キーを識別することによって始まる。異なるスーパークラスタ・キーを有するレコードは、構築によって、異なるクラスタに入るであろう。住所等の地理的なデータについては、郵便番号が適切なスーパークラスタ・キーであることが多い。整合する変形郵便番号を有するレコードは、当該スーパークラスタに含まれる場合がある。整合しない郵便番号を有するレコードは、高い確率で別のクラスタに属すると予想され、それ故に、性能を向上させるために、スーパークラスタ・キーを導入することによってクラスタを演算する際に、それらは排除される。

【 0 2 4 3 】

各スーパークラスタ内で、データは、会社名のフィールドの長さの降順且つ会社名の降順でソートされ、再現性のある順序で、最も長い名前を最初にクラスタ化アルゴリズムに提供する。スーパークラスタ・グループにおける最初のレコードは、最初のクラスタの主レコードとされる。その後の各レコード（ここでは、現レコードと呼ばれる）は、当該クラスタの主レコードの会社名に照らして現レコードの会社名を採点することによって、存在する各クラスタの主レコードと比較される。スコアが疑わしい整合閾値を超える場合、当該クラスタは、現レコードについての疑わしいクラスタのリストに加えられる。現レコードを存在する全ての主レコードと比較した後、疑わしいクラスタのリストが空である場合は、現レコードが、新しいクラスタの主レコードとされる。疑わしいリストが入力を 1 つだけ有しており、且つスコアが整合閾値を超えている場合は、現レコードは、当該疑わしいリスト上のクラスタに加えられる。疑わしいリストが 2 つ以上の入力を有している場合は、現レコード上の会社名が、当該疑わしいリスト上のクラスタの各々における全ての

10

20

30

40

50

レコードに照らして採点される。現レコードは、それが整合閾値を超える最高スコアを有するクラスタに加えられる。2つ以上のクラスタにおけるレコードについての最高スコアによる互角の整合が存在する場合、現レコードは、このようなクラスタの中の最初のクラスタに加えられる。整合閾値を超えるスコアが無い場合、現レコードは、新しいクラスタの主レコードとなる。

【0244】

このアルゴリズムは、2つの重要な特徴を有する。複数のクラスタへの曖昧な整合は最初に整合したクラスタの方に決定されるので、クラスタの中には、曖昧なメンバで比較的過剰に入力されているものもある。また、レコードがアルゴリズムに提供される順序は、特定のメンバシップ決定に影響を及ぼす。会社名の長さ及び値についての初期ソートは、名前の固定された順序を確立することによって、この問題を改善することを目的とする。下記において考察する部分的メンバシップの観念は、クラスタ・メンバシップの曖昧さをより正確に反映する、より高品位なソリューションを与える。

10

【0245】

曖昧なメンバシップの例を、以下の会社名の集合に示す。

【0246】

【表4】

ACME Services Australia Limited

ACME Services Canada Limited

ACME Services Limited

20

【0247】

特定の採点において、ACME Services Australia LimitedのACME Services Canada Limitedに対するスコアは、整合閾値0.75未満である0.65であり、これらのレコードは別個のクラスタに配置される。ACME Services Limitedは、両方のクラスタに対して等しいスコア0.95を有する。それは、ACME Services Australia Limitedと最初に出くわすので、ACME Services Australia Limitedのクラスタのメンバとなる。

【0248】

(2) 部分的メンバシップ

前のセクションにおける最初の例において、レコードA及びFは、クラスタC及びEの両方に属する。クラスタにおけるレコードの全ての出現が1とカウントされる場合、6つのレコードしか存在しないにもかかわらず、クラスタC及びEにおけるレコードの総カウントは8である(グループC(C、A、B、D、F)に5、及びグループE(E、A、F)に3)。この場合、部分的メンバシップを使用して、総カウントを維持することができる。1つのデータが2つ以上のグループに属する場合、そのデータの出現は1未満、即ち分数としてカウントされる。しかしながら、そのデータの全ての出現の合計は、やはり1となっており、総カウントを保つべきである。

30

【0249】

アレンジによっては、例えば、上述の曖昧さの尺度を使用して、ある要素が特定のグループに属する確度を反映するように、そのグループにおける当該要素の部分的メンバシップを定義してもよい。

40

【0250】

例えば、レコードAが40%の確率でグループCに属し、60%の確率でグループEに属すると仮定する。グループCにおけるレコードAには0.4の部分的メンバシップを割り付けることができ、グループEにおけるレコードAには0.6の部分的メンバシップを割り付けることができる。

【0251】

同様に、レコードFが10%の確率でグループCに属し、90%の確率でグループEに

50

属すると仮定する。グループCにおけるレコードFには0.1の部分的メンバシップを割り付けることができ、グループEにおけるレコードFには0.9の部分的メンバシップを割り付けることができる。

【0252】

レコードA及びFに割り付けられた部分的メンバシップにより、総カウントは、グループCのカウント(1+1+1+0.1+0.4=3.5)とグループEのカウント(1+0.9+0.6=2.5)との合計(6となる)である。故に、総カウントは維持される。

【0253】

部分的メンバシップの発端は、特定の要素のメンバシップに関する不確かさであり、各グループの総メンバシップは不確かさの度合い(即ち、誤差の範囲)によってのみ知られる。各グループの総カウントは、誤差の範囲によって調整される、全体の及び部分的な、メンバシップの合計として表現することができる。この範囲は、メンバシップについての全ての不確かな決定が包括又は排他の何れかに分類されるとみなすことによって得られる最大及び最小の境界値によって示すことができる。これらは、クラスタ間でのメンバの分布についての最悪の事態のシナリオに相当する。ここで、境界値について、Cにおける総メンバシップは3.5(3, 5)であろう。これは、Cのメンバの予想数は3.5であるけれども、Cは最少で3つのメンバ、最大で5つのメンバを有することを言うものと解釈される。同様に、Eにおける総メンバシップは2.5(1, 3)であろう。Eの予想される総メンバシップは2.5であるけれども、Eは最少で1つのメンバ、最大で3つのメンバを有する。

【0254】

異なるクラスタに属する境界値は相関性があるけれども、ここで使用される観念はそれを示唆しない。データセットにおける異なるレコード間の相関関係は可能であり、境界値を演算する際に考慮に入れられるべきである。例えば、時として、A及びFが、どのクラスタに属するかを知ること無く、それらが同じクラスタに無いことを知ることは可能である。

【0255】

上記第2の例において、特に、2つ以上のクラスタに対する曖昧な整合又は疑わしい整合が存在する場合に、上記疑わしい閾値を超える主レコードと整合する全てのレコードを各クラスタに関連付けることによって、クラスタ化の品質を高めることができる。当該整合の品質は、疑わしい各レコードに照らして記録され、曖昧さの尺度によって定量化されるべきである。部分的メンバシップは、あるアレンジにおいて、完全な(full)メンバとは別個に保持され、それらの部分的メンバシップの尺度によって表示されるであろう。例えば、あるクラスタのメンバは、(部分的メンバシップを有する完全なメンバと共に)部分的メンバシップの降順で列挙することができる。

【0256】

部分的メンバシップが一緒に特定されたレコードをリンクし且つ部分的メンバシップの割り当てを特定したルール、事象又は決定を識別するために使用することができるルール表示を各レコードにも添付すべきである。このルール・ラベルは、レコードを異なる部分的メンバシップと組み合わせる際に部分的メンバシップを調整する場合に有用であろう。

【0257】

ある見地からは、部分的メンバシップは、メンバシップの不確かさから生ずる曖昧さを反映するけれども、別の見地からは、部分的メンバシップは、2つの部門のために働いている従業員の例におけるように、単純に、複数のクラスタ間でのメンバシップの割り当てである。不確かな場合、既知知識における変化が、メンバシップの割り当てを変化させることが予想される。あるいは、部分的メンバシップを、単純に、信頼性のあるものとして許容してもよい。上記比率を現実のものとして許容するのに全くコストは掛からない。

【0258】

部分的メンバシップが、異なるクラスタに属する要素の確度を表現する場合、部分的メ

10

20

30

40

50

ンバシッは常に負でない値であり、異なるクラスタに属する要素の部分的メンバシッの合計は1であろう。

【0259】

しかしながら、部分的メンバシッは、アレンジによっては、負の値であってもよい。しかしながら、異なるクラスタに属するオブジェクトの部分的メンバシッの合計は、やはり1であるように制約されなければならない。

【0260】

アレンジによっては、要素の部分的メンバシッは、当該要素と上記主レコードとの間の距離の関数として、又は当該要素と上記主レコードとの間の整合スコアの関数として、定義することができる。ファジーなスコアから部分的メンバシッを構築する1つの方法は、上述のように、曖昧さの尺度によるものである。異なるファジーなスコアは、要素と上記主レコードとの間の異なる距離を反映し、故に、曖昧さの異なる尺度を反映する。ファジーなスコアは変形と上記種レコードとの間の類似性を反映し、確率と同じではないことが多いということに留意されたい。

【0261】

(h) ファジーなデータ操作

(1) 部分的メンバシッの存在下でのフィルタリング

選択基準を適用して、共通の性質を共有するレコードの部分集合を分離することが有用であることが多い。例えば、国際的なレコードのデータセットにおいて、特定の国からのレコードを選択することができる。当該選択操作（時として、「フィルタリング」と称される）は、選択を特定する表現において使用されるフィールドがキーである必要は無いので、キーに基づくものとはみなされない。レコードが複数のクラスタにおいて部分的メンバシッを有することが許される場合、フィルタリングは、部分的メンバの幾つかを失わせる場合がある。その結果は、1未満であるかもしれない選択された部分集合に亘るレコードに関連付けられた総メンバシッ割り当てである。これについての説明は、総割り当ては、選ばれた部分集合におけるメンバシッを、当該選ばれた部分集合の外にある選択肢に照らして測定するということである。

【0262】

ACME Services Limitedは、ACME Services Australia Limitedを含むグループに対して0.5の割り当てと、ACME Services Canada Limitedを含むグループに対して0.5の割り当てとを有する。データセット全体に亘るACME Services Limitedについての総割り当ては1.0である。フィルタを適用してCanadaと関連するレコードのみを保持する場合、ACME Services Limitedは、結果として得られるデータセットにおいて、0.5の総割り当てを有するであろう。このことは、ACME Services Limitedは、Canada部分集合中に無い選択肢については50%の可能性を有するのに対して、Canada部分集合中にある選択肢については50%の可能性を有することを示している。

【0263】

(2) キーによる並列処理及び部分的メンバシッ

並列処理において、キーの値に基づいて、異なる処理パーティションにレコードを割り当てることができる（時として、「キーによる分割」と称される）。レコードが曖昧なクラスタ・メンバシッを有することが許される場合、各クラスタに関連付けられたキーに基づいて分割を行ってもよい。この分割スキーム下では、所定のレコードに関連付けられたパーティション内での総割り当てが1未満であってもよい。このことの解釈は、フィルタリングについてのそれと似ている。それは、そのパーティションには無い選択肢に照らして、そのパーティションへのレコードの割り付けを測定する。

【0264】

ACME Services Limitedは、ACME Services Australia Limitedを含むグループに対して0.5の割り当てと、ACME

10

20

30

40

50

Services Canada Limitedを含むグループに対して0.5の割り当てとを有すると仮定する。キーによる分割操作は、ACME Services Australia Limitedを含むグループを1つのパーティションに、ACME Services Canada Limitedを含むグループをもう1つのパーティションに割り当てることができる。ACME Services Limitedのレコードに関連付けられた後者のパーティションにおける総割り当ては、何等かの他のパーティション（複数であってもよい）それがある0.5の選択肢に対して、0.5である（ACME Services Canada Limitedのクラスタとのその関連性を反映している）。

【0265】

10

よく知られているデータ操作の並列バージョンは、パーティション間でのやりとりが全く無い、単一のパーティション内でのそれらの挙動によって定義することができる。あるレコードについてのパーティション内での総割り当てが1未満である場合、このことは、ここで定義された意味に解釈される。

【0266】

(3) ロールアップ及び部分的メンバシップ

ロールアップ操作は、個々のレコードのレベルからグループのレベルへとデータを統合又は要約する。厳密なキーの場合、キーグループは、共通キー（値）を共有するレコードの集合として定義される。クラスタの場合、グループは、1つ以上のレコードが1つ以上のグループのメンバである可能性を伴う、メンバが比較によって特定されるレコードの集合として定義される。

20

【0267】

クラスタ・グループにおける付加的な（また、例えば、対数を加えることによる、乗法的な）数学的統合（時として、「演算加法的測度（computing additive measures）」と呼ばれる）は、重み付けについての割り当て尺度を使用する加重統合として行われる。境界値は、（集合への部分的割り付けを有する全てのレコードが当該集合に含まれているか又は当該集合から排除されているかの何れかである）選択肢についての（重み付けをしない）統合を演算することによって演算される。以下のレコードのリストは、会社名に基づくクラスタ・グループである。

【0268】

30

【表5】

クラスタ・キー	割り当て尺度	会社名	カウント
cl	1.0	ACME Services Australia Limited	80
cl	1.0	ACME Services (AUS) Limited	60
cl	0.5	ACME Services Limited	100

【0269】

クラスタにおける総カウントを特定するロールアップは、重み付き合計が $80 * 1.0 + 60 * 1.0 + 100 * 0.5 = 190$ であり、境界値が $80 * 1.0 + 60 * 1.0 + 100 * 0.0 = 140$ （排他的）及び $80 * 1.0 + 60 * 1.0 + 100 * 1.0 = 240$ （包括的）である。クラスタ・グループにおける総カウントについての結果は、190（140, 240）と表現することができる。

40

【0270】

部分的割り付けを有するレコードが集合に含まれるか又は集合から排除されるかの何れかである極端な場合を考察することによって、非付加的な要約を行う。割り付け尺度を使用して、部分的メンバを含めることによって得られる結果に信頼度を割り付けることができることが多い。例として、何等かの二次的なキーに基づいて、クラスタ・グループ内でレコードをソートしてもよく、そのソート順序において、どのレコードが最初であるのかをロールアップによって特定してもよい。以下のリストは、レコードの前記リストをカウ

50

ントの降順でソートしたものである。

【 0 2 7 1 】

【表 6】

クラスタ・キー	割り当て尺度	会社名	カウント
c1	0.5	ACME Services Limited	100
c1	1.0	ACME Services Australia Limited	80
c1	1.0	ACME Services (AUS) Limited	60

【 0 2 7 2 】

このソート順序でのクラスタ・グループにおける最初のレコード（即ち、最大カウント）を特定するロールアップにより、以下の境界の結果が与えられる。

【 0 2 7 3 】

【表 7】

(包括的) c1 0.5 ACME Services Limited 100

(排他的) c1 1.0 ACME Services Australia Limited 80

【 0 2 7 4 】

0.5 の割り当て尺度を包括的な結果に規定して、包括的結果に関連付けられる信頼度を示すことができる。この例においては、排他的結果は最悪の事態の結果とみなすことができる（最大値が少なくとも 80 である）。

【 0 2 7 5 】

部分的メンバシップの存在下でのロールアップ操作は、並列して実行してもよい。このことを理解するために、先ず付加的な数学的統合について考察する。これは加重和である。このような合計は、別個に演算される部分和に分け、次いで組み合わせることができる。各部分和は、その独自の並列パーティションにおいて演算することができる加重和である。これは、付加的ロールアップの並列化である。

【 0 2 7 6 】

定義によって全ての部分的メンバが排除されるので、付加的ロールアップ及び非付加的ロールアップの両方の排他的境界の演算は並列化可能である。従って、当該演算により、（殆どの場合において）並列化可能である通常のロールアップが減少する。

【 0 2 7 7 】

包括的境界の演算は、包括の二重カウントを防ぐ特定の条件下で並列化可能である。クラスタ内でのロールアップについて、各部分的メンバは当該クラスタ内で 1 回だけ発生する。故に、包括的境界は、何れのメンバも二重にカウントすること無くクラスタ内での部分和の合計として演算することができる。

【 0 2 7 8 】

ロールアップをクラスタに亘って実行する場合、異なるクラスタにおいて出現する同じレコードからの寄与があるかもしれない。このことは、当該レコードの各インスタントと関連付けられた重み付けが加わり、新しい全体としての重み付けが与えられるので、問題無い。しかしながら、包括的境界については、各レコードは 1 回だけ含まれるようにしなければならない。一般に、このことは、レコードが発生する経過を何等かの方法で追うことを必要とし、この操作は並列化可能ではない。

【 0 2 7 9 】

しかしながら、個々のレコードが、クラスタ化に先立って、キー（例えば、rec_key）によって識別され、（クラスタ化後の）データが rec_key に基づいて並列分割される場合、同じ rec_key を有する全てのレコードは同じパーティションにおいて発生するであろう。このパーティション内でのロールアップは、たとえクラスタに亘って行われても、全ての関係のあるレコードが存在するので、包括的境界を正しく演算することができる。個々のレコードにおいて、2 つ以上のパーティション上にインスタンスを

10

20

30

40

50

有するものは無く、従って、二重カウントの可能性は無いので、次いで、パーティションを亘る包括的境界を、パーティションを亘って安全に組み合わせることができる。

【0280】

(4) 検索

厳密なデータ操作においては、厳密なキーが使用される。例えば、検索操作において、キーが使用され、そのキーに厳密に整合する全てのレコードが取り出される。ファジーなデータ操作においては、ファジーなキーが使用される。

【0281】

アレンジによっては、例としての検索操作を使用して以下に説明するように、一連の厳密なデータ操作として、ファジーな操作が実行される。

【0282】

上記顧客住所データベースの例において、会社は、John Smithと名付けられた人物に送られた全ての保留中の郵便小包を見付け出すことに興味を覚える。その目的に、検索操作を使用することができる。検索操作は、二値キー「John Smith; ox2 6qt」を使用して実行することができる。ここで、郵便番号「ox2 6qt」は、John Smithと関連付けられた正しい郵便番号である。しかしながら、この厳密な検索操作は、入力を行う際に郵便局員によってなされた誤植のために、John Smithに送られた保留中の郵便小包を取り出さず、「John Smith; ox2 6qx」又は「John Smith; ox2 6qt」というキーと誤って関連付けられた保留中の郵便小包を取り出すであろう。

【0283】

この制約を克服するために、ファジーなキーと関連付けられる検索である、ファジー検索を使用することができる。ファジーなキーは、主キー及びそのキーの指定された距離内に収まる全ての変形を含むキーの集合の1つである。上記顧客住所データベースにおいては、上記保留中の郵便小包の検索において使用されるファジーなキーを、主キーである「John Smith; ox2 6qt」に加えて、2操作の距離内に収まる全ての変形を含むように定義することができる。

【0284】

アレンジによっては、ファジーなキー（主キーであるJohn Smith ox2 6qt及び4つの変形であるox2 6qx、ox2 6qy、ox2 6qy及びox2 6qt）に基づいて実行される上記ファジー検索操作は、以下のやり方で実行することができる。ステップ1において、主キーである「John Smith; ox2 6qt」に基づく厳密な検索を実行する。次いで、ステップ2から5において、上記ファジーなキーの一部である変形に基づいて厳密な検索を4回実行する。最後のステップ（ステップ6）において、上記ステップ1から5から取り出された結果を組み合わせる。当該組み合わせられた結果は、ファジーなキーを使用するファジー検索の結果である。

【0285】

会社の名前又は住所等の多単語フィールドにおいて検索を行う場合、前の手順に先だって、そのフィールドの変形の集合を特定し、直接的に使用することが可能ではない場合がある。2つの代替戦略を用いることができる。2 Plater Drive, OxfordのACME Service Ltdと関連付けられる全ての会計レコードについての検索がなされていると仮定する。第1の戦略においては、単一単語フィールド（例えば、郵便番号ox2 6qt）を検索キーとして使用する。当該キー「ox2 6qt」又は固定された距離以内のその変形の1つを有する全てのレコードが、「見込みのある整合」又は「有力候補」として、厳密な検索によって取り出される。各有力候補の会社の名前及び住所は、問合せ用の会社の名前及び住所に照らして、別個に採点される。

【0286】

概して、住所は、最初に、当該住所を構成する全てのフィールドを単一の文字列に連結することによって比較される。これには幾つかの利点がある。先ず、それは、異なるソースからの住所を単一の共通フォーマットに整える。句採点関数は単語の欠落や単語の順序

10

20

30

40

50

の変化に寛容であるので、ソースとなる住所の要素を保持している元のフィールドが不十分に又は一貫性無く埋められていても問題無い。同様に、住所を構文解析して標準的な住所の要素とする必要無しに、共通フォーマットに達する。構文解析は、概して、ポスタル・アドレス・ファイル等の正確な参照用データを必要とするけれども、特に外国においては、ポスタル・アドレス・ファイルは入手可能ではない場合がある。また、構文解析は、計算的には、相対的に高価であるので、それを回避することが、住所の比較の実行をより良好にする。連結された住所フィールドの比較が要領を得ない場合には、採点を改良するためのバックアップ的な試みとして、構文解析を使用してもよい。

【0287】

結果として得られた会社の名前及び住所のスコアは、(ユーザ指定の)整合閾値及び疑わしい閾値と比較される。整合閾値を超えるスコアは、比較された2つの句が整合を構成するのに十分に類似していることを示す。疑わしい閾値を超えるスコアは、それらの句が類似しているけれども、信頼度を以て整合を特定するのに十分な程には近くないことを示す。レコード全体が整合であるか否かについてのスコアは、個々のフィールド又はフィールドの組み合わせ(即ち、連結された住所)についてのスコアを組み合わせることによって得られる。どの情報がどれほど近く一致していなければならないかを特定する基準は、ユーザが指定する。例えば、市のフィールドは、郵便番号が一致する場合は、一致しなくても許される場合がある。(英国においては、例えば、何が有効な住所を構成するのかにおいて、驚くべき寛容度がある。町の名前及び番地の有無の両方についての変動が許容される。)

【0288】

第2の戦略は、1つ以上のフィールドから単語を選択し、これらをファジー検索用キーワードとして使用することである。このような単語は、フィールドにおける主要な単語として、又はそれらの有意性に基づいて、選ぶことができる。有意性は、フィールドにおける単語又はその変形の1つの出現数のそのフィールド(複数であってもよい)が埋められた(入力された)回数に対する比率のネガティブ・ログから演算される。

【0289】

ファジー検索用キーワードが選ばれると、各単語(及びその変形)がソース・レコードにおいて単語毎に調べられる。これは、選ばれた単語を所定のフィールドに含むソース・レコードのインデックスのリストを返す。検索用キーワードが、ポスタル・アドレス・ファイルの何れかの住所フィールドにおいて探される「Lower」、「Islip」であると仮定する。「Lower」という単語は、組織、通りの名前及び町のフィールドにおいて出現するかもしれない。「Islip」は町としてのみ出現するかもしれない。各単語についての検索により、上記単語が見付かるレコードのインデックスのリストが与えられる。これら2つのリストの積集合を得る(intersect)ことにより、両方の単語(又はそれらの変形)を含むレコードの集合が与えられる。これらは有力候補であるレコードの集合を形成する。有力候補であるレコードは特定の数の検索用の単語を含むことが知られているので、それらは効率的に事前承認される。共通して保持されることが要求される検索用の単語が多いほど、フィールド間のスコアがより高くなると見込まれる。

【0290】

例えば、2つ以上の検索用単語を含むレコードのみを保つためにフィルタを適用した後、全部のレコードが取り出され、互いに照らして採点される。結果として得られるスコアは、降順にソートされる。曖昧さの尺度が整合の集合について演算され、これにより、参照に対する問合せの整合の品質を特徴付ける情報が増える。

【0291】

(5) ファジーな結合

レコードを取り出すのにルックアップを使用する代わりに、参照用データセット全体が読み込まれ、問合せデータセットに照らして処理される結合が行われることを除き、ファジーな結合はファジー検索に類似している。このことは、性能及び制御の両方について有用である場合がある。参照用データセットが、通常のルックアップとしてメモリに納める

10

20

30

40

50

には大き過ぎる場合、代わりに、ディスク上に保持されている（ことによると、ブロック圧縮された）ロード可能なルックアップとしてアクセスしてもよい（Ab Initioのロード可能なルックアップを参照）。各検索語が処理される際に、ルックアップ・テーブルの適切なページがディスクからアクセスされる。参照用データセットの十分に大きい断片（例えば、10%）にアクセスする必要がある場合は、ランダム・アクセス検索を始めるのではなく、単一のパスにおいてソートされた順序で参照用データセット全体を読み込む方が、結局、より効率的となる。これにより、検索プロセスによって繰り返し実行されるディスク・アクセスが減少する。

【0292】

制御の見地からは、検索が複数の参照用データセットへのアクセスを必要とする場合、結合が、これらの参照用データセットからのデータを採点のために一緒にするのに、より便利な方法である。問合せが単一の入力における顧客の名前及び住所を問合せが含むけれども、参照用データは顧客の名前及び住所を別個のテーブルに保持していると仮定する。顧客の名前と住所とを繋ぐ第3のリンク・テーブルが存在する場合がある。対象となる参照用データが異なるテーブルに保持されている場合、インデックスは異なるデータセットにおけるレコードを参照するので、検索用インデックスについてインデックス・ファイルを検索することによって返されるレコード・インデックスを直接的に比較することは可能ではない。問合せにキーが割り付けられている場合は、各データセットに照らして別個の検索を行うことができ、これらの検索の結果を、問合せキーに基づく結合によって組み合わせ、有力候補を取ってきて、採点することが可能となる。

10

20

【0293】

曖昧さの尺度がファジーな結合の出力の一部として演算される場合、各々が部分的メンバシップを有する複数の整合を結合操作の結果として生じさせることが可能である。例えば、番地が無い住所をポータル・アドレス・ファイルに照らして結合して、郵便の宛先のファイルがその住所に記録されている組織の名前を取得すると仮定する。ポータル・アドレス・ファイルにおいて適切な通りにある3つの住所には組織が無く、故に当該組織を識別するために、これらの住所を組み合わせることができる。他の整合する住所には、2つの別個の組織であるACME Ltd及びStandard Corp.が存在する。曖昧さの均等分配尺度により、同等な存在の数がカウントされる。当該結合の出力は、最初は、整合する住所の各々に1つずつ、5つのレコードであろう。続いて、（曖昧さの尺度を特定するために）当該住所における組織に対するロールアップを行うことにより、当該組織が、空欄（3/5）、ACME Ltd（1/5）、Standard Corp.（1/5）のいずれかである確度を示す曖昧さの均等分配尺度が得られる。次いで、この結果を、各々が別個の組織及び関連付けられた部分的メンバシップを伴う、3つのレコードに正規化する。

30

【0294】

【表 8】

問合せ:住所に基づく結合により組織を取得する

問合せ住所: Lower Street, ox2 6qt

組織	住所
--	2 Lower Street, ox2 6qt
--	3 Lower Street, ox2 6qt
ACME Ltd	4 Loower St, ox2 6qt
--	5 Lower St, ox2 6qy
Standard Corp.	8 Lower St., ox2 6qt

10

結果:

組織	住所	割り当てキー	部分的メンバシップ°
--	Lower Street, ox2 6qt	a1	0.6
ACME Ltd	Lower Street, ox2 6qt	a1	0.2
Standard Corp.	Lower Street, ox2 6qt	a1	0.2

【 0 2 9 5 】

20

(6) ソート及び部分的メンバシップ

メンバシップが部分的である場合にレコードを順序付ける(ファジーな)ソートは定義が容易である。曖昧な関連性を有するレコードについては、その割り当て(曖昧さの尺度)と一緒に、各選択肢についてのレコードが作り出される。前の例における参照用レコードが、問合せレコードに照らして、結合の結果でソートされると仮定する。ルールは、部分的メンバシップは、完全なメンバシップ後に、メンバシップの降順でソートされるというものである。他のフィールドについての下位ソート(subsort)は、部分的メンバシップが適用された後に適用される。従って、1つのフィールドに対する部分的メンバシップは、後のキーに対して上位にある。これにより、追加のソート・フィールドの適用は、より高いレベルにある順序を変更すること無く、確率された順序でレコードをソート

30

【 0 2 9 6 】

【表 9】

ソート{組織部分住所}:

--	2 Lower Street, ox2 6qt	--	1.0
--	3 Lower Street, ox2 6qt	--	1.0
--	5 Lower St, ox2 6qy	--	1.0
--	Lower Street, ox2 6qt	a1	0.6
ACME Ltd	4 Loower St, ox2 6qt		
ACME Ltd	Lower Street, ox2 6qt	a1	0.2
Standard Corp.	8 Lower St., ox2 6qt		
Standard Corp.	Lower Street, ox2 6qt	a1	0.2

40

【 0 2 9 7 】

この種のファジーなソートを適用した場合、組織への(ソートされた)ファジーなロールアップは、全てのデータが見られるまで一時的な結果を格納することを必要とせず、組織名に基づいて行うことができる。これは、ソートされたデータの主な用途の1つであ

50

り、ロールアップ操作は、各キー・グループが完成する際に完了する。

【0298】

ファジーなマージ操作は、ファジーなソートに似ている。それは、単に、その入力上の各レコードに順序付けルールを適用して、そのソート順助において、どのレコードが次にあるのかを特定する。上記データセットが以下のデータセットとマージされると仮定する。

【0299】

【表10】

ACME Ltd	Lower Street, ox2 6qt	a2	0.9	
Standard Corp.	Lower Street , ox2 6qt	a2	0.1	10

【0300】

マージされたデータは以下の通りである。

【0301】

【表11】

マージ{組織部分住所}:

--	2 Lower Street, ox2 6qt	--	1.0	
--	3 Lower Street, ox2 6qt	--	1.0	
--	5 Lower St, ox2 6qy	--	1.0	20
--	Lower Street, ox2 6qt	a1	0.6	
ACME Ltd	4 Lower St, ox2 6qt			
ACME Ltd	Lower Street, ox2 6qt	a2	0.9	
ACME Ltd	Lower Street, ox2 6qt	a1	0.2	
Standard Corp.	8 Lower St., ox2 6qt			
Standard Corp.	Lower Street, ox2 6qt	a1	0.2	
Standard Corp.	Lower Street , ox2 6qt	a2	0.1	30

【0302】

(h) ファジーなデータ操作の有用性：間違い及び不確かさに取り組み、会計の完全性を保つ。

【0303】

上記ファジーな検索の例において示したように、ファジーな検索操作により、厳密なキーを使用する従来の検索では見逃すであろうレコード（例えば、誤植を含むレコード）を取り出すことができる。

【0304】

また、上述のように、データの分類が保留中の結果に依存する場合、クラスタ化又は部分的メンバシップを使用して、不確かさを正確に捉えることができる。直感的には、クラスタ化又は部分的メンバシップは、別個のファジーではない操作の組み合わせ又は連なりと同等であると見ることができる。しかしながら、クラスタ化又は部分的メンバシップは、より良い取り扱い及び予測を可能とするであろう。組織の年間予算が、当該組織が非慈善団体であるか否かに関する判決に依存している上記の例において、有利又は不利な判決が出される確率に基づいて年間予算を準備することができる。

【0305】

より具体的には、年間予算は以下のように予定しておくことができる。

【0306】

10

20

30

40

【数 1 4】

<p>予定される納税額＝</p> <p>非慈善団体の地位に基づく納税額× 不利な判決の確率</p> <p>+</p> <p>慈善団体の地位に基づく納税額× 有利な判決の確率</p>
--

10

【0307】

上記の式を使用して納税予想額を計算することは、組織の財政状態のより良好な全体像を提供し、経営幹部によるリスク評価を容易にする。それはまた、より信頼性の高い数字を下流のアプリケーションに与え、例えば、組織の財務見通しをより良好に予想することを可能とする。

【0308】

部分的メンバシップはまた、会計の完全性を保つのに有用である。例えば、ABC社の人材データベースにおいて、マーケティング部門及びR&D部門におけるJohn Smithのメンバシップがそれぞれ0.5である場合、彼の医療費が二重にカウントされるようなことはないということが判る。

20

【0309】

上述のアプローチは、コンピュータ上で実行されるソフトウェアを使用して実装することができる。例えば、かかるソフトウェアは、少なくとも1つのプロセッサ、少なくとも1つのデータ記憶システム（揮発性及び不揮発性メモリ、並びに/又は記憶素子）、少なくとも1つの入力装置又は入力ポート、及び少なくとも1つの出力装置又は出力ポートを各々が含む、1つ以上のプログラムされたまたはプログラム可能なコンピュータ・システム（分散型、クライアント/サーバ型、又はグリッド型等の種々のアーキテクチャのものであってよい）上で実行される、1つ以上のコンピュータ・プログラムにおけるプロシージャを形成する。上記ソフトウェアは、例えば、演算グラフの設計及び構成に関する他のサービスを提供する、より大きいプログラムの1つ以上のモジュールを形成していてもよい。上記グラフのノード及び要素は、コンピュータ可読媒体に格納されたデータ構造、又はデータ・リポジトリに格納されたデータ・モデルに合致する他の系統的データとして、実装することができる。

30

【0310】

上記ソフトウェアは、汎用の又は特殊用途のプログラム可能なコンピュータによって読み取ることができる記憶媒体（例えば、CD-ROM）上に提供したり、又は実行時に、ネットワークの通信媒体越しに（伝搬信号中にコード化して）コンピュータに配信したりすることができる。上記機能の全ては、特殊用途のコンピュータ上で、又はコプロセッサ等の特殊用途のハードウェアを使用して、行うことができる。上記ソフトウェアは、当該ソフトウェアによって規定される計算の異なる部分が異なるコンピュータによって行われる分散方式にて実装することもできる。このような各コンピュータ・プログラムは、汎用の又は特殊用途のプログラム可能なコンピュータによって読み取ることができる記憶媒体又は記憶装置（例えば、固体メモリ若しくは媒体、又は磁気媒体若しくは光媒体）上に格納したり、又はダウンロードしたりして、コンピュータ・システムが上記記憶媒体又は記憶装置を読み取って本明細書に記載されたプロシージャを実行しようとする際に、上記コンピュータが構成され、運転されるようにするのが好ましい。また、本発明のシステムは、コンピュータ・プログラムによって構成される、コンピュータ可読記憶媒体として実装されると考えることもできる。ここで、このように構成された記憶媒体は、コンピュータ

40

50

・システムを、特定の予め定められた方式で運転して、本明細書に記載された機能を実行する。

【0311】

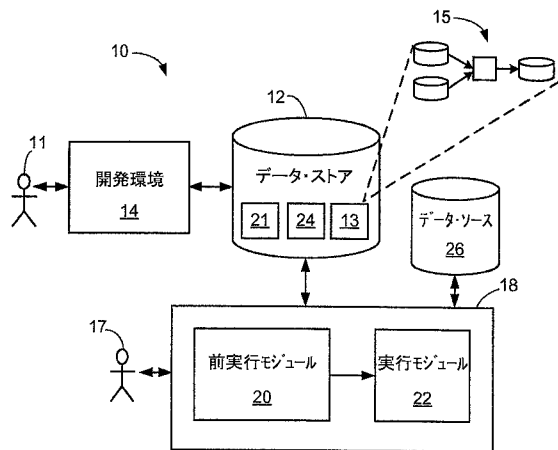
本発明の幾つかの実施態様について説明してきた。しかしながら、本発明の精神及び範囲から逸脱すること無く、様々な変形を行うことができることが理解されるであろう。例えば、上述のステップの幾つかは、順序に依存しないものであってもよく、従って、上述のものとは異なる順序で実行することができる。

【0312】

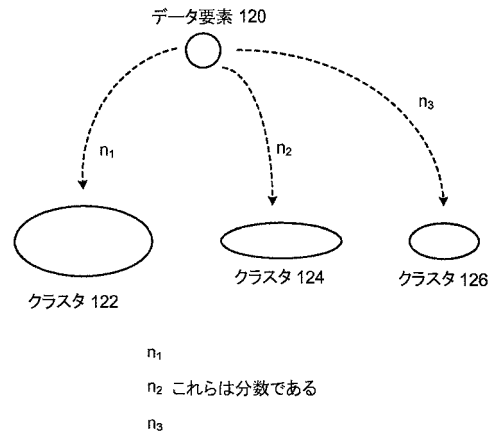
先の記述は本発明の範囲を説明することを目的とするものであり、添付した特許請求の範囲によって定義される本発明の範囲を限定することを目的とするものではないこともまた理解されるべきである。例えば、上述の機能ステップの多くは、全体としての処理に実質上の影響を及ぼすこと無く、異なる順序で実行することができる。他の態様は、以下の特許請求の範囲の範囲内にある。

10

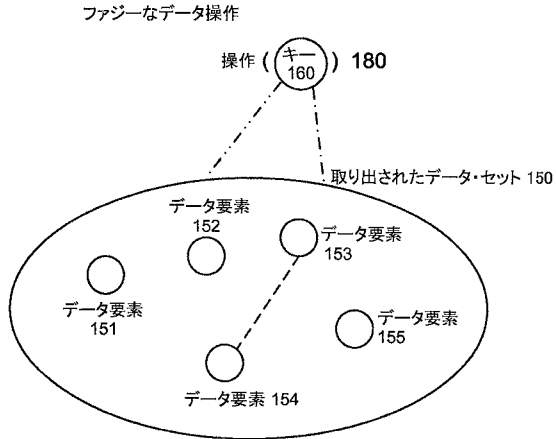
【図1】



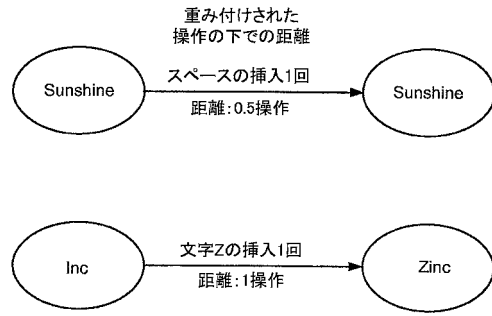
【図2A】



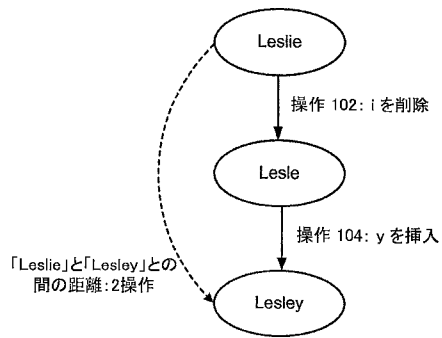
【図 2 B】



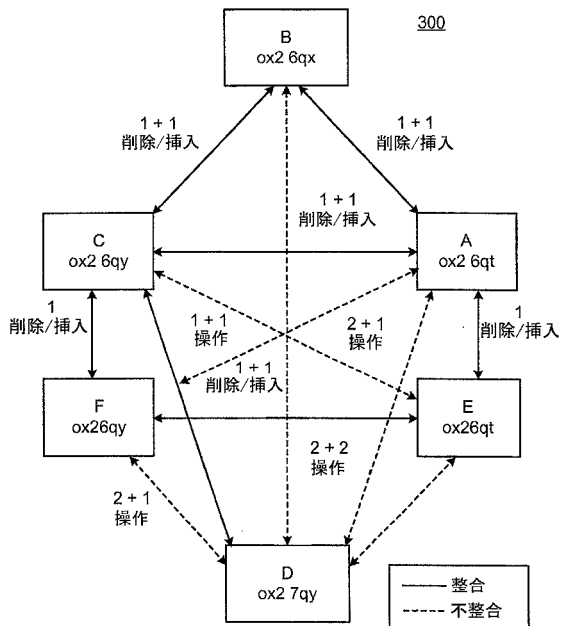
【図 2 D】



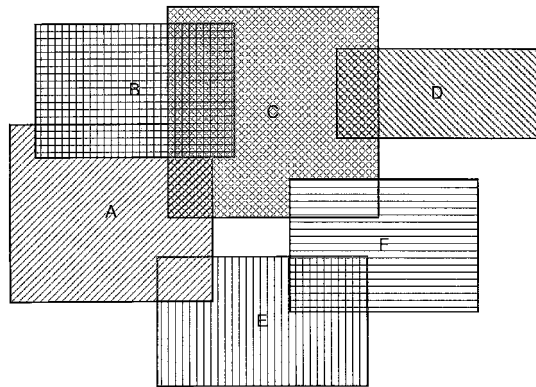
【図 2 C】



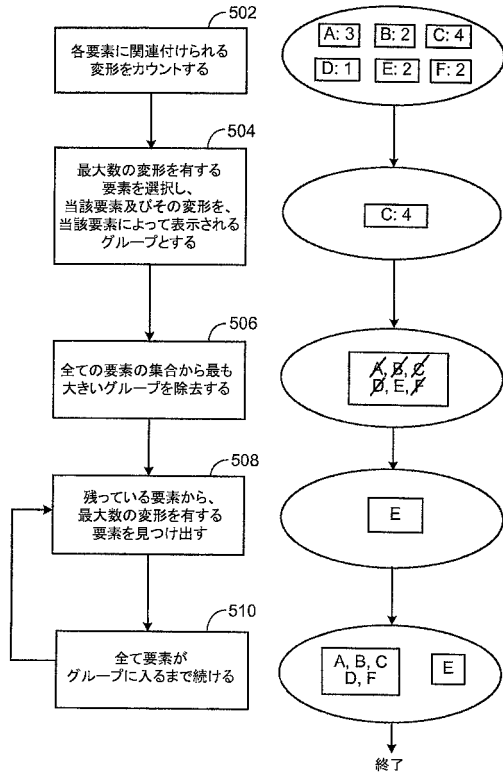
【図 3】



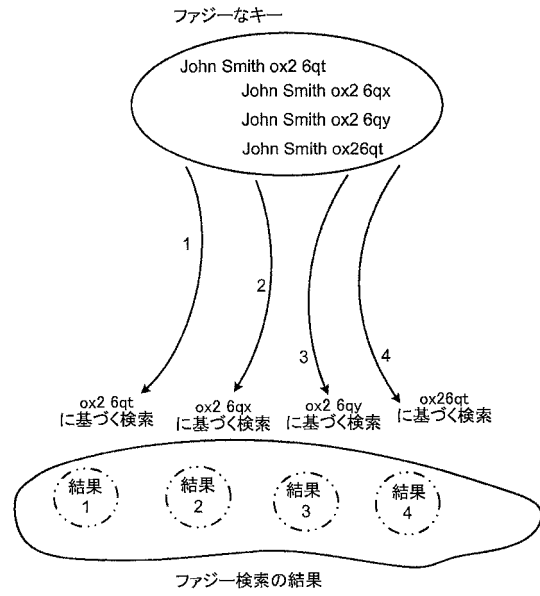
【図 4】



【図5】



【図6】



フロントページの続き

(56)参考文献 特開2003-006226(JP,A)
特開平09-044518(JP,A)
特開平06-044309(JP,A)

(58)調査した分野(Int.Cl., DB名)

G06N 3/00-99/00
G06F 17/30