

(19) 日本国特許庁(JP)

(12) 特 許 公 報(B2)

(11) 特許番号

特許第5394245号
(P5394245)

(45) 発行日 平成26年1月22日(2014. 1. 22)

(24) 登録日 平成25年10月25日(2013. 10. 25)

(51) Int.Cl.

F I

G O 6 F 17/30 (2006.01)

G O 6 F 17/30 3 5 O C

G O 6 F 17/30 1 7 O B

請求項の数 16 (全 14 頁)

(21) 出願番号	特願2009-533936 (P2009-533936)	(73) 特許権者	509116772
(86) (22) 出願日	平成19年10月23日(2007. 10. 23)		モンロ、ドナルド・マーティン
(65) 公表番号	特表2010-507856 (P2010-507856A)		イギリス国、ビーエー１１・６アールエス
(43) 公表日	平成22年3月11日(2010. 3. 11)		、サマーセット、ベッキントン、グース・
(86) 国際出願番号	PCT/GB2007/004035		ストリート、ザ・レイズ 6
(87) 国際公開番号	W02008/050107	(74) 代理人	100108855
(87) 国際公開日	平成20年5月2日(2008. 5. 2)		弁理士 蔵田 昌俊
審査請求日	平成22年10月25日(2010. 10. 25)	(74) 代理人	100091351
(31) 優先権主張番号	11/585, 358		弁理士 河野 哲
(32) 優先日	平成18年10月23日(2006. 10. 23)	(74) 代理人	100088683
(33) 優先権主張国	米国 (US)		弁理士 中村 誠
		(74) 代理人	100109830
			弁理士 福原 淑弘
		(74) 代理人	100075672
			弁理士 峰 隆司

最終頁に続く

(54) 【発明の名称】 ファジーデータベースマッチング

(57) 【特許請求の範囲】

【請求項 1】

サンプルレコードと、データベースの複数の記憶されているレコードとの間で可能性ある一致を識別する方法において、

前記記憶されているレコードのそれぞれから複数のインデックスコードを抽出し、前記インデックスコードは、インデックスコード空間内に含まれていることと、

サンプルレコードからサンプルインデックスコードを抽出することとを含み、

前記方法は、

前記インデックスコード空間を規定するルックアップテーブルを維持し、前記ルックアップテーブルは、複数の行を持ち、それぞれの行は前記インデックスコード空間内の一意的なインデックスコードに対応していることと、

複数のレコード出現リストを維持し、前記レコード出現リストのそれぞれは、特定のインデックスコードに対応している前記ルックアップテーブル中の特定の行からリンクされており、前記レコード出現リストのそれぞれは、そこから前記特定のインデックスコードと、前記インデックスコード空間内で前記特定のインデックスコードに対する規定された近接性内にあるインデックスコードとが抽出された、記憶されているレコードを識別するものであることと、

前記サンプルインデックスコードをインデックスとして使用して、ルックアップテーブルをアドレス指定し、前記サンプルインデックスコードに関係する、対応する複数のレ

10

20

コード出現リストをルックアップすることと、

インデックスコードが抽出されるときに、前記記憶されているレコードによる一致を記録するヒストグラムを構築することと、

前記ヒストグラムからの前記レコード出現リスト内で識別されたそれぞれの記憶されているレコードの出現回数をカウントすることと、

所定の記憶されているレコードに対する前記ヒストグラムからの前記カウントが、要求されるしきい値を超える場合、前記サンプルとの可能性ある一致であるとして、前記所定の記憶されているレコードを識別することにより特徴付けられる方法。

【請求項 2】

前記規定された近接性は、規定されたハミング距離である、請求項 1 記載の方法。

【請求項 3】

前記規定されたハミング距離は、ユーザ選択可能である、請求項 2 記載の方法。

【請求項 4】

前記要求されるしきい値は、数値的しきい値である、請求項 1 記載の方法。

【請求項 5】

前記要求されるしきい値は、記憶されているレコード毎の、レコード出現リストの平均数の関数である、請求項 1 記載の方法。

【請求項 6】

前記複数のインデックスコードは、前記複数の記憶されているレコードから抽出された、前記インデックスコード空間内のすべてのインデックスコードを規定する、請求項 1 記載の方法。

【請求項 7】

前記複数のインデックスコードは、あるサンプルレコードによって表示することができる前記インデックスコード空間内のすべての可能性あるインデックスコードを規定する、請求項 1 記載の方法。

【請求項 8】

前記複数のインデックスコードは、前記記憶されているレコードに対して、ハッシュのような演算を適用することによって生成される、請求項 1 記載の方法。

【請求項 9】

複数の規定された近接性を確立することと、それぞれのインデックスコードと近接性との組み合わせに対して、独立したレコード出現リストを維持することを含む、請求項 1 記載の方法。

【請求項 10】

前記識別ステップは、ユーザ選択された規定された近接性に関連するリストを使用する、請求項 1 記載の方法。

【請求項 11】

前記サンプルレコードと、前記可能性ある一致のそれぞれとの間の関係をさらに解析する追加のステップを含む、請求項 1 記載の方法。

【請求項 12】

前記識別ステップは、複数の並列プロセッサの間で分けられ、それぞれが関係付けの結果をコンソリデータに送り、前記コンソリデータは、前記関係付けの結果に依存して、可能性ある一致として、記憶されているレコードを識別する、請求項 1 記載の方法。

【請求項 13】

サンプルレコードと、データベースの複数の記憶されているレコードとの間の可能性ある一致を識別するシステムにおいて、

前記記憶されているレコードから抽出された複数のインデックスコードを含むデータベースに結合されたコンピュータプロセッサを具備し、

前記システムは、前記サンプルレコードからサンプルインデックスコードを抽出するように構成されており、

10

20

30

40

50

前記システムは、

前記インデックスコード空間を規定するルックアップテーブルと、
複数のレコード出現リストと

により特徴付けられ、

前記ルックアップテーブルは、複数の行を持ち、それぞれの行は前記インデックスコード空間内の一意的なインデックスコードに対応しており、

前記レコード出現リストのそれぞれは、特定のインデックスコードに対応している前記ルックアップテーブル中の特定の行からリンクされており、前記レコード出現リストのそれぞれは、そこから前記特定のインデックスコードと、前記インデックスコード空間内で前記特定のインデックスコードに対する規定された近接性内にあるインデックスコードとが抽出された、記憶されているレコードを識別するものであることと、

10

前記システムは、

前記サンプルインデックスコードをインデックスとして使用して、ルックアップテーブルをアドレス指定し、前記サンプルインデックスコードに関係する、対応する複数のレコード出現リストをルックアップするようにと、

インデックスコードが抽出されるときに、前記記憶されているレコードによる一致を記録するヒストグラムを構築するようにと、

前記ヒストグラムからの前記レコード出現リスト内で識別されたそれぞれの記憶されているレコードの出現回数をカウントするようにと、

(c) 所定の記憶されているレコードに対する前記ヒストグラムからの前記カウントが、しきい値を超える場合、前記サンプルレコードとの可能性ある一致であるとして、前記所定の記憶されているレコードを識別するプロセッサと
を具備するシステム。

20

【請求項 14】

前記コンピュータプロセッサは、サンプルレコードからサンプルインデックスコードを抽出する第 1 のプロセッサと、前記サンプルレコードとの可能性ある一致であるとして、所定の記憶されているレコードを識別する第 2 のプロセッサとを備え、請求項 13 記載のシステム。

【請求項 15】

前記第 1 のプロセッサは、前記第 2 のプロセッサから離れている、請求項 14 記載のシステム。

30

【請求項 16】

前記第 1 のプロセッサは、複数の並列プロセッサを備え、それぞれが関係付けの結果をコンソリデータに送り、前記コンソリデータは、前記関係付けられた結果に依存して、可能性ある一致として、記憶されているレコードを識別する、請求項 13 記載のシステム。

【発明の詳細な説明】

【発明の分野】

【0001】

本発明は、データベースシステムの分野に関連する。より詳細には、本発明は、候補レコードをデータベース内のレコードに対して高い信頼度でファジーに一致させる、速度を改善する方法およびシステムに関連する。

40

【先行技術】

【0002】

さまざまな分野において、特定のサンプルレコードが、大規模データベース内に既に存在するか否かを、非常にすばやく決定できること、そして、もし存在する場合、1 つ以上の一致を識別することに対する増加している要望がある。1 つの特定の分野は、バイオメトリクスであり、バイオメトリクスでの要求は、特定のバイオメトリックサンプルを提供した個人が、既にデータベース中に存在するか否かを決定することである。

【0003】

説明したタイプのデータベースは、きわめて大規模であってもよく、サンプルレコード

50

と、データベース内のレコードのあらゆるものとの間で、完全一致解析を試みることは非実用的であるかもしれない。計算作業負荷を減少させるために、さまざまな事前選別プロセスが使用されているが、これらは、マッチングアルゴリズムの特色、または、マッチングされることになるデータの特色に依存することが多いので、これらの多くのものは、非常に制限された適用分野しか持たない。

【 0 0 0 4 】

他の応用においても起こることであるが、バイオメトリックデータのマッチングで特に起こる問題は、バイオメトリック測定値が、それらの本質によって、正確に再現されないことが多いことである。例えば、ある特定の個人の虹彩から繰り返し導出される、バイオメトリック測定値は、とりわけ、まぶたとまつ毛による虹彩のオクルージョンが画像ごとに異なることになるので、多少異なっていることが多い。結果として、バイオメトリックマッチングは、通常、完全一致というよりは、近似、または、“ファジー”一致の概念に依拠する。

【 0 0 0 5 】

一般的なシナリオは、特定の個人が、個々人の大規模データベース内に存在するか否かを決定する要望である。例えば、我々は、ある個人の虹彩スキャンを持っていたとしてもよく、国家安全保障データベースが、同一の個人の1つ以上の虹彩スキャンを既に含んでいるか否かについて知りたいかもしれない。サンプルの虹彩スキャンと、記憶されている虹彩スキャンとは、あらゆる観点で、同一でないことが多いので、必要とされる“ファジー”一致を達成する1つの方法は、ある領域にわたって検索することである。サンプルと記憶されているレコードとの両方をコードに変換した後、何らかの予め規定されたプロトコルにしたがって、我々は、ある記憶されたレコードと、我々がサンプルコードに十分近いとして考える、ある領域内の任意のコードとの間の一致を見つけることを試みることができる。代わりに、我々は、サンプルコードと、記憶されているコードのうちの1つに十分近い、ある検索領域内の任意のコードとの間の一致を見つけようとしてもよい。何れのケースでも、ファジー一致を実行するときに、コードの領域全体にわたって検索する必要性は、マッチングプロセスをかなり遅延させるかもしれない。

【 発明の概要 】

【 0 0 0 6 】

本発明の第1の観点にしたがうと、サンプルレコードと、複数の記憶されているレコードとの間で可能性ある一致を識別する方法が提供され、方法は、

(a) 記憶されているレコードから複数の特性を抽出し、前記特性は、特性空間内に含まれることと、

(b) それぞれの前記特性に対する、記憶されているレコードのレコード出現リストを維持し、記憶されているレコードのレコード出現リストから、前記特性と、前記特性空間内の前記特性に対する規定された近接内の特性とが抽出されていることと、

(c) サンプルレコードから特性を抽出することと、

(d) 所定の記憶されているレコードが、要求される数のレコード出現リスト中に現れる場合、サンプルとの可能性ある一致であるとして、所定の記憶されているレコードを識別することを含む。

【 0 0 0 7 】

本発明のさらなる観点にしたがうと、特性空間内の複数の特性を使用して、サンプルレコードと、複数の記憶されているレコードとの間の可能性ある一致を識別するシステムにおいて、システムは、

(a) それぞれの特性に対する、記憶されているレコードのレコード出現リストから、前記特性と、前記特性空間内の前記特性に対する規定された近接内の特性とが抽出されている、記憶されているレコードのレコード出現リストと、

(b) サンプルレコードから特性を抽出するプロセッサと、

(c) 所定の記憶されているレコードが、要求される数のレコード出現リスト中に出現する場合、サンプルとの可能性ある一致であるとして、所定の記憶されているレコードを識

10

20

30

40

50

別するプロセッサとを具備する。

【0008】

このような方法は、データベース内に新しいレコードを登録するときに任意の追加の労力をかけて、非常に高速な候補マッチングを提供する。新しいレコードの登録の頻度に比べて、マッチングがより頻繁に行われるとき、このトレードオフは、十分に価値がある。

【0009】

いくつかの実施形態では、サンプルレコードに対する特性のマッチングのために、また、可能性ある一致として、記憶されているレコードを識別するために、独立したプロセッサを使用してもよい。これらのプロセッサは、独立したコンピュータであってもよく、互いに離れていてもよい。

10

【0010】

ある特定の例では、記憶されているレコードの完全なコレクションを含むメインデータリストが、特性リストとは別のものとして保持されてもよい。このことは、ローカルプロセッサが、ローカルに取得された虹彩スキャンのようなサンプルレコード上に初期解析を実行することを可能にする。一度、可能性ある一致のリストが識別されると、次に、このリストは、遠隔サーバに送られてもよく、遠隔サーバにおいて、サンプルと、可能性ある一致のそれぞれの完全な、エンコードされた虹彩スキャンとを比較することによって、より詳細な解析が実行される。

【0011】

このアプローチは、システムの設計者が、エンコードされた虹彩スキャンの全体のデータベースの完全なコピーを多数のユーザに配信する必要がないというさらなる利点を持っている。代わりに、それぞれのユーザは、特性のリストを単に受け取るだけであり、ローカルに実行されることになる初期解析に対しては、このような特性のリストで十分である。1つ以上の可能性ある一致が見つかり、システムが、集中ロケーションに自動的にレポートしてもよく、集中ロケーションにおいて、完全なレコードに対して、さらなる解析が実行されてもよい。

20

【0012】

本発明は、さまざまな方法で実行されてもよく、いくつかの特定の実施形態を、添付の図面を参照して、例として、ここで説明する。

【図面の簡単な説明】

30

【0013】

【図1】図1は、本発明の実施形態にしたがった、データベース構造を示す。

【図2】図2は、マッチングプロセスを例示するヒストグラムである。

【図3】図3は、別の例示的なヒストグラムである。

【図4】図4は、いくつかの例示的なハードウェアを示す。

【発明の詳細な説明】

【0014】

以下の詳細な説明では、多数の特定の詳細を述べて、特許請求されている主題の完全な理解をもたらす。しかしながら、特許請求されている主題は、これらの特定の詳細なしで実行されてもよいことが当業者によって理解されるだろう。他の例では、よく知られた方法、手続、構成部品、および/または、回路は、詳細に説明しない。

40

【0015】

以下の詳細な説明のいくつかの部分は、コンピュータ、および/または、コンピューティングシステムメモリ内のような、コンピューティングシステム内に記憶される、アルゴリズム、ならびに/あるいは、データビットおよび/またはバイナリデジタル信号上の動作の象徴的表現に関して提示した。これらのアルゴリズム的記述および/または表現は、データ処理技術における当業者によって使用され、彼らの作業の内容を他の当業者に伝える。アルゴリズムは、ここで、および一般的に、演算の筋の通ったシーケンス、および/または、望ましい結果へと導く類似の処理であるとして考えられる。動作および/または処理は、物理的な量の物理的な操作を含んでもよい。一般的に、必ずしも必要ではないが

50

、これらの量は、記憶され、転送され、結合され、比較され、および／またはそうでなければ、操作されることができる、電子のおよび／または磁気的信号の形態をとってもよい。主に共通の利用のために、時として、これらの信号をビット、データ、値、エレメント、シンボル、キャラクタ、ターム、数字、数値、および／または類似物として呼ぶことが便利であることが分かっている。しかしながら、これらのすべてのものと、類似のタームとは、適切な物理的量に関係付けられることになり、単に便利なラベルにすぎないことを理解すべきである。そうではないと特に明記しない限り、以下の説明から明らかなように、この明細書の説明全体にわたって、“処理”、“演算”、“計算”、“決定”、および／または、類似物のようなタームを利用することは、コンピュータまたは類似の電子のコンピューティングデバイスのような、コンピューティングプラットフォームの動作および／または処理を指し、コンピューティングプラットフォームは、コンピューティングプラットフォームのプロセッサ、メモリ、レジスタ、ならびに／あるいは、他の情報記憶装置、送信、および／または表示デバイス内の物理的電子のおよび／または磁気的な量、ならびに／あるいは、他の物理的量として表されるデータを操作し、および／または、変換する。

10

【0016】

簡潔さのために、以下の説明は、バイオメトリック分野の例示的な実施形態に向けられることとする。説明することになる実施形態では、特定の個人から虹彩スキャンが取られ、同一の個人の虹彩スキャンが、国家安全保障データベースのような大規模なデータベース内に既に存在しているか否かを決定することが必要である。

20

【0017】

もちろん、この特定の例は、本発明の背景にある、一般的な原則を説明するために、単に使用されており、同じ技法を他の分野中でも等しく適用できることになっていくことが理解されるだろう。その最も広い形態における本発明は、データベース内に保持される、任意の特定のクラスまたはタイプのデータに限定されることはなく、使用されるマッチングアルゴリズムの詳細に限定されることはない。

【0018】

例示的な実施形態のデータベース構造を、図1に概念的に示した。特定の個人の詳細が、ケースリストまたはテーブル16内に保持されており、ケースリストまたはテーブル16のそれぞれの行17は、特定の個人の虹彩スキャンを表す。理想的には、それぞれの個人が、単一の虹彩スキャンによって表現されることになるが、もちろん、一般的な国家安全保障データベースでは、実際には、少なくとも数人の個人の複数のスキャンがあることになる。それぞれの行または虹彩スキャンレコードは、列18、20、22を含み、これは、それぞれ、システム内で使用するための一意的な虹彩スキャン参照番号と、知られている場合は、個人の名前と、国家安全保障または社会安全保障コードのような外部識別子とをそれぞれ保持する。

30

【0019】

それぞれのレコードに対する完全な虹彩スキャンが、独立したデータリストまたはテーブル10内で保持され、独立したデータリストまたはテーブル10のそれぞれの行11は、個人のスキャンを表す。このテーブルは、2つの列からなり、最初の列12は、前述のような一意的な参照番号であり、第2の列14は、特定の個人の虹彩スキャン、または他のバイオメトリックレコードを一意的に識別する、何らかのエンコードされた表現を保持するとして考えられてもよい。

40

【0020】

それぞれの登録されたケース（虹彩スキャン）は、複数の属性、特性、またはコードにしたがって分類され、これらは、未処理の虹彩スキャン、または、より一般的にはエンコードされたスキャンデータ14のいずれかから、抽出または導出される。

【0021】

コードは、スキャンの人的識別可能な特徴を表現していてもよいが、必ずしもそうでなくてもよい。例えば、コードのいくつかのものは、眼の色を表現していてもよく、他のも

50

のは、虹彩内部の色素の量と、輝度変化のような特性を表現していてもよい。代わりに、エンコードされたスキャン 14 は、データストリームに何らかの関数が単に適用された単純な結果としてのコードを備える、純粋なデータストリームとして取り扱われてもよい。前述のハッシュ関数以外の、さらなる可能性は、データストリーム内の特定のグループのビットの存在または、不在を検索することである。いずれにせよ、それぞれの個人のレコード 11 から、複数のコードが一般的に抽出されることが理解されるだろう。

【0022】

これらのコードをインデックスとして使用するのを容易にするために、(以下でより詳細に説明することになるように、)コードは一般的に数値であるとして制約され、特定の予め規定された範囲にあることが多い。バイオメトリクス応用では、例えば、コードは、16 バイナリビットによって規定されていてもよく、65536 個の可能性あるコードが起こることを可能にしている。好ましくは、所望のコードだけが利用可能であるように、未処理の、または、エンコードされたデータからこれらのコードを生成させる関数または演算は、これらの可能性ある出力範囲において制限されていてもよい。代わりに、出力の実際の範囲が、所望の範囲内の数値コードのリストに再マッピングされてもよい。要求される場合には、(表示していない)マッピングテーブルを使用してもよい。説明している例では、利用可能なコード P_n は、0 から 65535 の範囲内の整数であり、それぞれが 16 ビットコードとして記憶されていることが仮定されるだろう。したがって、フル 16 ビットまで、 $P_1 = 001$ 、 $P_2 = 010$ 、 $P_3 = 011$ 、 $P_4 = 100$ 、 $P_5 = 101$ 等である。

【0023】

その中でそれらが現れるケース(虹彩スキャン)にしたがって、コードを分類するために、複数のリストまたはテーブル 28 が、可能性ある 65536 個のコードのそれぞれに対して 1 つ、維持される。簡潔さのために、これらのリストの 5 つのものだけを図 1 に示した。理解されるように、コード値 1 に対するリスト 40 は、単一の行だけを含み、この行は、名前 A だけがこのコードを生成したことを示す。コード値 2 を表現するリスト 41 は、この例では、登録された虹彩スキャンのどれもそのコードを生成していないので、何のデータも含んでいない。リスト 42、43 は、それぞれコード値 3 と 4 を表現しており、単一のスキャンだけに関連する。表 44 は、名前 A と B に対する虹彩スキャンがそれぞれコード値 5 を生成することを示す。

【0024】

本質的ではないが、それぞれのテーブルまたはリスト 28 が、それぞれの行 29 中に、該当するコードに対応する単一のレコードに対する、一意的な参照 18 を単に含むことが一般的に好ましい。

【0025】

リスト 28 に加えて、第 2 の系列のリスト 30 が維持され、これらのリストは、単に個別のコードに関連しているのではなく、ハミング距離のような何らかの所望のメトリックにしたがった、対応しているベースコードからの所定の距離であるコードに関連する。

【0026】

図 1 によって与えられる例を参照すると、コード 1 から 5 の間のハミング距離は、表 1 で与えられる。ハミング距離は、2 つのコードの間で異なっているビットの数である。例えば、コード 1 と 2 との間のハミング距離は、1 のためのコード(001)と、2 のためのコード(010)との間で 2 ビットが変更されているので、2 である。

【表 1】

コード	1 (001)	2 (010)	3 (011)	4 (100)	5 (101)
1	0	2	1	2	1
2	2	0	1	2	3
3	1	1	0	3	2
4	2	2	3	0	1
5	1	3	2	1	0

表1 数1から5に対するバイナリコードの間のハミング距離

【0027】

示した例において、テーブル30は、テーブル28の対応するベースコードから、ちょうど1のハミング距離を持っているコードを解決するケースに関連するデータを含む。したがって、例えば、テーブル51は、P1(0001)からちょうどH=1距離であるすべてのコードに対するデータを含む。しかしながら、H=1のコードが、既にテーブル40中で出現している場合、これらは効率のために、テーブル51から省略される。図1の例では、コード1は、コード3および5から、H=1距離である。3に対するベースリストからの名前Dが適格となり、5に対するベースリストからの名前AとBが適格となる。しかしながら、名前Aは、コード1に対するベースリスト中で既に発生しているので、したがって、名前DとBだけが、テーブル51中に含まれる。同様に、テーブル52が、P2(0010)からちょうどH=1距離であるすべてのコードに対するデータを含む、等である。しかしながら、H=1のコードが、既にテーブル41中で出現している場合、これらは効率のために、テーブル52から省略される。例において、名前Dが3に対するベースリストからテーブル52に追加されるように、コード3は、コード2からH=1距離である。テーブル41が空であるので、テーブル52から省略すべきである、コード2からのH=1のコードは何もない。

【0028】

第3の系列のテーブル32は、対応するベースコードから、H=2距離であるコードを解決するケースの詳細を含む。要求される場合、(示していない)H=3、H=4等に対する、さらなる系列のテーブルもまた提供されてもよい。

【0029】

実施形態を図示するのにハミング距離を使用したことと、他の任意の便利なメトリックを使用してもよいことが理解されるだろう。取り扱う特定の応用にしたがって、要求されるメトリック(例えば、ハミング距離)を選択してもよく、さらに、固定されていても、ユーザ選択可能であってもよい。(示していない)より洗練された実施形態では、コードは、対応する多次元空間内で測定されている要求されるメトリック内の、多次元のものであってもよい。

【0030】

データベース内で、新しい虹彩スキャンが登録されることになる时候はいつでも、その詳細が、ケースリスト16と、データリスト10とに追加され、新しいスキャンに対する対応するコードが計算され、および/または、決定される。次に、スキャンの一意的な参照番号18が、個別のリスト28、30、32に対して、適切に追加される。望ましい場合、1つ以上の新しいコードがコードリスト24に追加されてもよく、この中で、個別のテーブル28が自動的に作成され、データベース内のそれぞれの虹彩スキャンをチ

10

20

30

40

50

チェックして、その参照番号を、新しく作成されたテーブルの1つ以上に対して追加する必要があるか否かを決定する。

【0031】

ここで、マッチングのタスク、または、言い換えると、ある未知の虹彩スキャンがデータベース内のスキャン14のうちの1つに一致するか否かを決定することに戻る。スキャンをエンコードされたデータに対してマッチングさせることは、計算的に冗長であるかもしれないので、スキャンをエンコードされたデータに対してマッチングさせるというよりはむしろ、代わりに、サンプルを進めて、それから1つ以上のコード値を抽出する。登録されているスキャンに元々適用されたのと同じ関数を適用することによって、1つ以上のサンプルコードが生成される（もちろん、この例におけるコードは、すべて、整数であり、範囲0から65536内にある）。 10

【0032】

どのスキャンが、それぞれのサンプルコードに対応するかを見つけるために、それぞれのコードnを、ルックアップテーブル25に対するインデックス70として使用し、このテーブルは、コード値1、2、3を保持する、メモリ中のそれぞれのエリアをポイントするポインタP1、P2、P3・・・を含む。それぞれのリストが、メモリ中で互いに引き続いて起こる、特定の公称コード値に集める場合、単一のポインタ（プラスオフセット）だけが要求されるだろう。代わりに、独立したポインタが、系列28、系列30、および系列32内でそれぞれのリストに対して提供されてもよい。別の可能性は、リスト28のそれぞれが、30中の対応するリストをルックするポインタを持つこと等であるだろう。 20

【0033】

一度、適切なテーブルが識別されると、システムは、次に、特定のハミング距離のすべてのテーブルにわたって、それぞれのケースの出現回数のヒストグラムを構築することによって、一致候補を識別するために進む。図2は、サンプルスキャンが、コード値1と3を生成し、1までのハミング距離を使用して、一致候補が識別されることになる例を図示する。関連するH=1テーブル51、53とともに、コード1と3に対して、ベーステーブル40、42中のレコードを見ることによって、このことが達成される。ベースコードテーブルは、2つのヒット、すなわちAとDを生成し、一方、H=1テーブルは3つの追加ヒットA、B、およびDを生成する。 30

【0034】

カウントに対してしきい値が適用され、少なくともしきい値にスコア付けされる任意のレコードが、一致候補であるとして考えられる。ここで、しきい値を1ととる場合、一致候補は、スキャンA、B、およびDである。2のしきい値において、候補はAとDである。

【0035】

図3は、コード1と3を生成する、同一のサンプルに対するヒストグラムを示すが、今回は、2までのハミング距離に対して試みた。ここで、ベーステーブル40、42からのヒットは、AとDであり、H=1テーブル51、53からの追加のヒットは、A、B、およびDであり、H=2テーブル61、63からの追加のヒットは、BおよびCである。1のしきい値を適用することは、一致候補として、A、B、C、およびDを与え、他方、より高いしきい値2を適用することは、A、B、およびDを候補として戻す。 40

【0036】

図2、3中で、カウントをヒストグラムとして示したが、他のカウント方法も等しく使用されてもよく、いずれにせよ、実際のヒストグラムは一般的に作図されないだろう。

【0037】

（示していない）代替の実施形態では、第2の系列のテーブル30は、対応するベースコードから、ちょうどH=1距離であるコードだけからのデータのみならず、最大その距離までのすべてのコードからのデータを含む。このような配置では、それぞれのH=1テーブルは、対応するベーステーブルのすべてのデータを含み、それぞれのH=2テーブルは、対応するH=1テーブルのすべてのデータを含む等であるだろう。 50

【 0 0 3 8 】

しきい値およびノまたはハミング距離に対して、適切な値を選択することによって、応用にしたがって、システムの出力応答を調整してもよい。これらの値のいずれか、または、両方は、固定であってもよく、プログラミング可能に可変であってもよく、ユーザ可変のものであってもよい。任意のアプリケーションでは、ランタイムにおいて、これらのパラメータのいずれかまたは両方の適切な値をユーザが選択できることが便利であるかもしれない。

【 0 0 3 9 】

いくつかの応用では、より複雑なマッチングアルゴリズムが構想されてもよい。例えば、異なるハミング距離に対して、異なるしきい値を使用してもよい。システムは、さまざまなハミング距離において、候補を自動的に選択してもよく、異なる距離において、それぞれの選択を比較または組み合わせて、候補リストの改善された複合リストを生成してもよい。

10

【 0 0 4 0 】

全体のマッチングプロセスを高速化するために、事前選択プロセスが、多数のケースを考察から除外する必要の度合いにしたがって、必要である場合、しきい値およびノまたはハミング距離の選択が決定されてもよい。単純なカウントと、固定されたしきい値の使用は、一致を不一致から分ける便利な方法であるが、他のアルゴリズムを等しく使用することもできる。1つの可能性あるアプローチは、例えば、すべてのケースにわたって集められた平均特性カウントよりも、ある固定パーセンテージ分より高い特性カウントを持っている、すべてのケースを一致候補として選択することであるだろう。

20

【 0 0 4 1 】

評価されることになるサンプルのサイズに依拠して、そのサンプルを全体的に使用する必要はないかもしれない。例えば、サンプルが、ある書籍のいくつかの章からなる場合、1ページのテキストだけに基づいて、事前選択を実行することが十分であるかもしれない。

【 0 0 4 2 】

誤った棄却の危険性が受入可能なまでに低くなるように、ほとんどのアプリケーションにおいて、特性の選択、マッチング基準、および、解析されることになるサンプルのサイズが、選択されることになるだろう。

30

【 0 0 4 3 】

一度、一致候補のリストが選択されると、前述の手続の内の1つを使用して、何らかの便利なマッチングアルゴリズムを使用して、可能性あるもののそれぞれに対して、より詳細な一致を実行してもよい。説明したテキストの例では、サンプルスキャンが、データベース内の候補に対して、より洗練されているが、より遅いアルゴリズムを使用して、比較されてもよい。

【 0 0 4 4 】

1つの実施形態では、データベース自体は、事前の、およびノまたは最終のマッチングが発生するのと、同じコンピュータ、または同じロケーションに保持されていてもよい。代わりに、事前のマッチングがローカルコンピュータにおいて保持されているコードリストにしたがって実行され、事前の一致がリモートコンピュータに送られて、詳細なマッチングが発生するというように、プロセスが分散されていてもよい。このような配置は、(すべての記憶されているスキャンを表す完全なデータを含む)一次データリスト10が、中央ロケーションにおいて保持されることを可能にし、ここで、ローカルマシンは、個別のケース出現リスト28、30、32だけを保持すればよい。

40

【 0 0 4 5 】

図4に示した、別の実施形態では、本発明のプロセスは、並列で動作する、複数のコンピュータまたはプロセッサを使用することによって、さらに高速化されている。ユーザのコンピュータ32は、マッチングタスクを制御装置34に送り、制御装置34は、これを分けて、複数のコンピュータまたはプロセッサ36の間で分散させる。それぞれのプロセ

50

ッサ 36 は、特定のコード、または、グループのコードを取り扱うように命令されていてもよい。代わりに、制御装置 34 は、何らかの別の方法で、作業を分けてもよい。プロセッサ 36 は、それらの結果をコンソリデータ 38 に送り、コンソリデータ 38 は、(例えば、図 2、3 で図示した手続を使用して、) 可能性ある一致の選択を最終的に行う。次に、可能性あるもののリストが、要求されるように、詳細なマッチングを実行するコンピュータまたはプロセッサ 42 に送られ、あるいは、参照番号 40 で示すように、さらなる解析のためにユーザ 32 に戻される。

【0046】

もちろん、特定の実施形態を説明してきたが、特許請求されている主題は、特定の実施形態、または、実現に限定されていないことが理解されるだろう。例えば、1つの実施形態は、デバイスまたはデバイスの組み合わせ上で動作するように実現されるような、ハードウェア中のものであってもよく、他方、別の実施形態は、例えば、ソフトウェア中のものであってもよい。同様に、ある実施形態は、例えば、ファームウェア中で、または、ハードウェア、ソフトウェア、および/または、ファームウェアの任意の組み合わせとして実現されてもよい。同様に、特許請求されている主題は、この観点における範囲に制限されていないが、1つの実施形態は、記憶媒体のような1つ以上の物品を含んでもよい。例えば、1つ以上のCD-ROM、および/または、ディスクのような、この記憶媒体は、命令を記憶してもよく、命令は、コンピュータシステム、コンピューティングプラットフォーム、または、他のシステムのようなシステムによって実行されるときに、例えば、前述した実施形態のうちの1つのような、特許請求されている主題にしたがった方法の実施形態が実行されることに帰結してもよい。1つの可能な例として、コンピューティングプラットフォームは、1つ以上の処理ユニット、または、プロセッサ、ディスプレイ、キーボードおよび/またはマウスのような1つ以上の入出力デバイス、ならびに/あるいは、スタティックランダムアクセスメモリ、ダイナミックランダムアクセスメモリ、フラッシュメモリ、および/またはハードドライブのような1つ以上のメモリを含んでもよい。

【0047】

これまでの説明において、特許請求されている主題のさまざまな観点を説明した。特許請求されている主題の完全な理解をもたらすために、説明の目的で、特定の数、システム、および/または構成を述べた。しかしながら、特許請求されている主題は、これらの特定の詳細なしで実行されてもよいことが、この開示の利益を受ける当業者にとって明らかになるだろう。他の例では、特許請求されている主題をあいまいにしないように、よく知られた特徴を省略し、および/または簡潔化した。ある特徴をここで図解および/または説明したが、当業者にとって、さまざまな修正、置換、変更、および/または、均等物が、ここで明らかになるだろう。したがって、添付の特許請求の範囲は、このような修正および/または変更のすべてが、特許請求されている主題の本来の精神の範囲内におさまることを意図している。

以下に、本願出願の当初の特許請求の範囲に記載された発明を付記する。

[1] サンプルレコードと、複数の記憶されているレコードとの間で可能性ある一致を識別する方法において、

(a) 前記記憶されているレコードから複数の特性を抽出し、前記特性は、特性空間内に含まれていることと、

(b) それぞれの前記特性に対する、記憶されているレコードのレコード出現リストを維持し、記憶されているレコードのレコード出現リストから、前記特性と、前記特性空間内の前記特性に対する規定された近接性内の特性とが抽出されていることと、

(c) サンプルレコードから特性を抽出することと、

(d) 所定の記憶されているレコードが、要求される数のレコード出現リスト中に現れる場合、前記サンプルとの可能性ある一致であるとして、前記所定の記憶されているレコードを識別することと

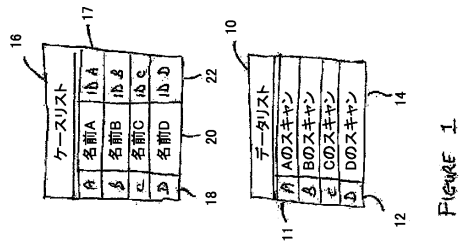
を含む方法。

[2] 前記規定された近接性は、規定されたハミング距離である、上記[1]の方法。

- [3] 前記規定されたハミング距離は、ユーザ選択可能である、上記 [2] の方法。
- [4] 前記要求される数は、数値的しきい値である、上記 [1] の方法。
- [5] 前記要求される数は、記憶されているレコード毎の、レコード出現リストの平均数の関数である、上記 [1] の方法。
- [6] 前記複数の特性は、前記複数の記憶されているレコードから抽出される前記特性空間内のすべての特性を規定する、上記 [1] の方法。
- [7] 前記複数の特性は、あるサンプルレコードによって表示することができる前記特性空間内のすべての可能性ある特性を規定する、上記 [1] の方法。
- [8] 前記複数の特性は、前記記憶されているレコードに対して、ハッシュのような演算を適用することによって生成される、上記 [1] の方法。 10
- [9] 前記サンプルレコードに演算を適用して、1つ以上のサンプル出力を生成することと、前記サンプル出力を使用して、ルックアップテーブルをアドレス指定することとを含み、前記ルックアップテーブル中のそれぞれの行はレコード出現リストをポイントする、上記 [1] の方法。
- [10] 特性が抽出されるときに、記憶されているレコードによる一致を記録するヒストグラムが構築され、前記ヒストグラムから、可能性ある一致として、レコードを識別する、上記 [1] の方法。
- [11] 複数の規定された近接性を確立することと、それぞれの特性と近接性との組み合わせに対して、独立したレコード出現リストを維持することとを含む、上記 [1] の方法 20
- [12] 前記識別ステップは、ユーザ選択された規定された近接性に関連するリストを使用する、上記 [11] の方法。
- [13] 前記サンプルレコードと、前記可能性ある一致のそれぞれとの間の関係をさらに解析する追加のステップを含む、上記 [1] の方法。
- [14] 前記識別ステップは、複数の並列プロセッサの間で分けられ、それぞれが関係付けの結果をコンソリデータに送り、前記コンソリデータは、前記関係付けの結果に依存して、可能性ある一致として、記憶されているレコードを識別する、上記 [1] の方法。
- [15] 特性空間内の複数の特性を使用して、サンプルレコードと、複数の記憶されているレコードとの間の可能性ある一致を識別するシステムにおいて、 30
- (a) それぞれの特性に対する、記憶されているレコードのレコード出現リストから、前記特性と、前記特性空間内の前記特性に対する規定された近接性内の特性とが抽出されている、前記記憶されているレコードのレコード出現リストと、
- (b) 前記サンプルレコードから特性を抽出するプロセッサと、
- (c) 所定の記憶されているレコードが、要求される数のレコード出現リスト中出现する場合、前記サンプルとの可能性ある一致であるとして、前記所定の記憶されているレコードを識別するプロセッサと
- を具備するシステム。
- [16] 前記抽出するプロセッサと、前記識別するプロセッサとは、共通のプロセッサからなる、上記 [15] のシステム。
- [17] 前記抽出するプロセッサは、前記識別するプロセッサから離れている、上記 [1 5] のシステム。 40
- [18] 前記抽出するプロセッサは、複数の並列プロセッサを備え、それぞれが関係付けの結果をコンソリデータに送り、前記コンソリデータは、前記関係付けられた結果に依存して、可能性ある一致として、記憶されているレコードを識別する、上記 [15] のシステム。

【図 1】

図 1



【図 2】

図 2

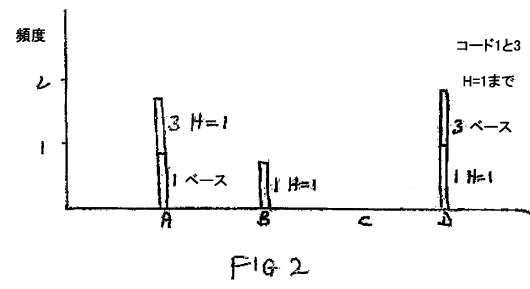
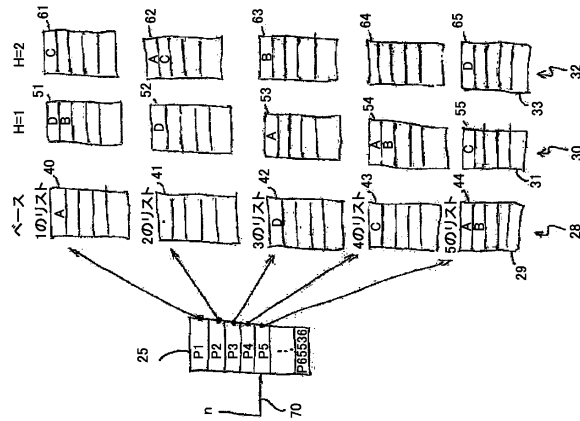


FIG 2



【図 3】

図 3

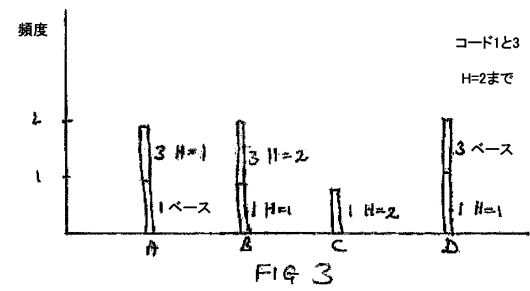


FIG 3

【図 4】

図 4

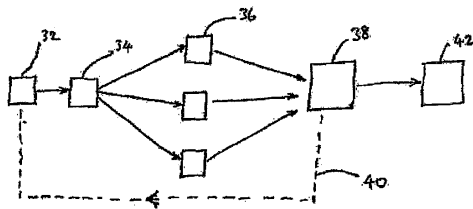


Figure 4

フロントページの続き

(74)代理人 100095441

弁理士 白根 俊郎

(74)代理人 100084618

弁理士 村松 貞男

(74)代理人 100103034

弁理士 野河 信久

(74)代理人 100140176

弁理士 砂川 克

(74)代理人 100100952

弁理士 風間 鉄也

(72)発明者 モンロ、ドナルド・マーティン

イギリス国、ビーエー１１・６アールエス、サマーセット、ベッキントン、グース・ストリート、
ザ・レイズ 6

審査官 佐藤 実

(56)参考文献 米国特許出願公開第２００３／００６１２３３（ＵＳ，Ａ１）

特表平０８－５０４９７９（ＪＰ，Ａ）

米国特許出願公開第２００６／０１０４４９３（ＵＳ，Ａ１）

米国特許出願公開第２００４／０２０２３５５（ＵＳ，Ａ１）

(58)調査した分野(Int.Cl.，ＤＢ名)

G 0 6 F 1 7 / 3 0