

(19) 日本国特許庁(JP)

(12) 特 許 公 報(B2)

(11) 特許番号

特許第6434162号  
(P6434162)

(45) 発行日 平成30年12月5日(2018.12.5)

(24) 登録日 平成30年11月16日(2018.11.16)

(51) Int.Cl.	F I
<b>G06F 17/30 (2006.01)</b>	G06F 17/30 350C
	G06F 17/30 414A
	G06F 17/30 412

請求項の数 16 (全 39 頁)

(21) 出願番号	特願2017-547261 (P2017-547261)	(73) 特許権者	000003078 株式会社東芝 東京都港区芝浦一丁目1番1号
(86) (22) 出願日	平成27年10月28日(2015.10.28)	(73) 特許権者	301063496 東芝デジタルソリューションズ株式会社 神奈川県川崎市幸区堀川町72番地34
(86) 国際出願番号	PCT/JP2015/080446	(74) 代理人	110002147 特許業務法人酒井国際特許事務所
(87) 国際公開番号	W02017/072890	(72) 発明者	浜田 伸一郎 東京都港区芝浦一丁目1番1号 株式会社東芝内
(87) 国際公開日	平成29年5月4日(2017.5.4)	(72) 発明者	小野 聡一郎 東京都港区芝浦一丁目1番1号 株式会社東芝内
審査請求日	平成30年1月26日(2018.1.26)		

最終頁に続く

(54) 【発明の名称】 データ管理システム、データ管理方法およびプログラム

(57) 【特許請求の範囲】

【請求項1】

蓄積するデータの特徴ベクトルである事例ベクトルに類似する周辺ベクトルを生成し、生成した前記周辺ベクトルに対応する前記事例ベクトルを特定するための索引情報を構築する索引構築部と、

任意の特徴ベクトルであるクエリベクトルを指定した検索要求に応じて、前記索引情報を用いて、前記クエリベクトルと完全一致する前記周辺ベクトルに対応する前記事例ベクトルを特定し、特定した前記事例ベクトルに基づく検索結果を出力する検索部と、を備え

る。  
前記索引構築部は、少なくとも、前記周辺ベクトルを格納する第1列と、該周辺ベクトルに対応する前記事例ベクトルに関する情報を格納する第2列とを列要素を持つテーブルと、該テーブルにおける前記第1列に対する索引とを含む前記索引情報を構築し、

前記検索部は、前記索引を用いて、前記クエリベクトルと完全一致する前記周辺ベクトルに対応する前記テーブルのレコードを求め、求めたレコードの前記第2列に格納された情報に基づいて前記事例ベクトルを特定する、データ管理システム。

【請求項2】

前記テーブルのデータ構造として、前記第1列に格納される前記周辺ベクトルをキーとし、前記第2列に格納される情報を値とする連想配列または連続メモリ配置型配列を用いる、請求項1に記載のデータ管理システム。

【請求項3】

前記索引構築部は、前記第 1 列および前記第 2 列に加えてさらに、前記周辺ベクトルの前記事例ベクトルに対する類似度を格納する第 3 列を列要素に持つ前記テーブルと、該テーブルにおける前記第 1 列および前記第 3 列に対する複合索引とを含む前記索引情報を構築し、

前記検索部は、前記複合索引を用いて、前記クエリベクトルと完全一致する前記周辺ベクトルであって、前記類似度が条件を満たす前記周辺ベクトルに対応する前記テーブルのレコードを求め、求めたレコードの前記第 2 列に格納された情報に基づいて前記事例ベクトルを特定する、請求項 1 に記載のデータ管理システム。

【請求項 4】

前記テーブルのデータ構造として、前記第 1 列に格納される前記周辺ベクトルおよび前記第 3 列に格納される前記類似度をキーとし、前記第 2 列に格納される情報を値とする連想配列または連続メモリ配置型配列を用いる、請求項 3 に記載のデータ管理システム。

【請求項 5】

蓄積するデータの特徴ベクトルである事例ベクトルに類似する周辺ベクトルを生成し、生成した前記周辺ベクトルに対応する前記事例ベクトルを特定するための索引情報を構築する索引構築部と、

任意の特徴ベクトルであるクエリベクトルを指定した検索要求に応じて、前記索引情報を用いて、前記クエリベクトルと完全一致する前記周辺ベクトルに対応する前記事例ベクトルを特定し、特定した前記事例ベクトルに基づく検索結果を出力する検索部と、を備え

、  
前記索引構築部は、前記周辺ベクトルを格納する第 1 列と、該周辺ベクトルの前記事例ベクトルに対する類似度を格納する第 2 列とを列要素に持つ第 1 テーブルと、該第 1 テーブルのレコードの行 ID を格納する第 1 列と、該レコードの前記周辺ベクトルに対応する前記事例ベクトルに関する情報を格納する第 2 列とを列要素に持つ第 2 テーブルと、前記第 1 テーブルにおける前記第 1 列および前記第 2 列に対する複合索引とを含む前記索引情報を構築し、

前記検索部は、前記複合索引を用いて、前記クエリベクトルと完全一致する前記周辺ベクトルであって、前記類似度が条件を満たす前記周辺ベクトルに対応する前記第 1 テーブルのレコードの行 ID を求め、求めた行 ID が格納された前記第 2 テーブルのレコードの第 2 列に格納された情報に基づいて前記事例ベクトルを特定する、データ管理システム。

【請求項 6】

蓄積するデータの特徴ベクトルである事例ベクトルに類似する周辺ベクトルを生成し、生成した前記周辺ベクトルに対応する前記事例ベクトルを特定するための索引情報を構築する索引構築部と、

任意の特徴ベクトルであるクエリベクトルを指定した検索要求に応じて、前記索引情報を用いて、前記クエリベクトルと完全一致する前記周辺ベクトルに対応する前記事例ベクトルを特定し、特定した前記事例ベクトルに基づく検索結果を出力する検索部と、を備え

、  
前記索引構築部は、前記周辺ベクトルの値と該周辺ベクトルの前記事例ベクトルに対する類似度の条件とに従って、前記周辺ベクトルに対応する前記事例ベクトルに関する情報を探索する複合索引を前記索引情報として構築し、

前記検索部は、前記複合索引を用いて、前記クエリベクトルと完全一致する前記周辺ベクトルであって、前記類似度が条件を満たす前記周辺ベクトルに対応する前記事例ベクトルを特定する、データ管理システム。

【請求項 7】

蓄積するデータの特徴ベクトルである事例ベクトルに類似する周辺ベクトルを生成し、生成した前記周辺ベクトルに対応する前記事例ベクトルを特定するための索引情報を構築する索引構築部と、

任意の特徴ベクトルであるクエリベクトルを指定した検索要求に応じて、前記索引情報を用いて、前記クエリベクトルと完全一致する前記周辺ベクトルに対応する前記事例ベク

10

20

30

40

50

トルを特定し、特定した前記事例ベクトルに基づく検索結果を出力する検索部と、を備え

前記検索部は、前記検索要求が出力件数の指定を含む場合に、前記事例ベクトルに対する前記周辺ベクトルの類似度の条件を厳しい方から段階的に変化させながら、前記クエリベクトルと完全一致する前記周辺ベクトルに対応する前記事例ベクトルを特定する処理を、特定した前記事例ベクトルの総数が指定された前記出力件数以上になるまで繰り返し、特定した前記事例ベクトルの総数が指定された前記出力件数以上になると前記処理を停止して、特定した前記事例ベクトルに基づく前記出力件数に近い件数の検索結果を出力する、データ管理システム。

【請求項 8】

蓄積するデータの特徴ベクトルである事例ベクトルに類似する周辺ベクトルを生成し、生成した前記周辺ベクトルに対応する前記事例ベクトルを特定するための索引情報を構築する索引構築部と、

任意の特徴ベクトルであるクエリベクトルを指定した検索要求に応じて、前記索引情報を用いて、前記クエリベクトルと完全一致する前記周辺ベクトルに対応する前記事例ベクトルを特定し、特定した前記事例ベクトルに基づく検索結果を出力する検索部と、を備え

前記検索部は、前記索引情報を用いた検索方式と、元の特徴ベクトルまたは縮約ベクトル空間に写像された特徴ベクトルを用いて線形検索する検索方式と、を含む複数の検索方式のうち、前記事例ベクトルを含むデータテーブルのレコード数に応じて選択された検索方式により前記事例ベクトルの検索を行う、データ管理システム。

【請求項 9】

蓄積するデータの特徴ベクトルである事例ベクトルに類似する周辺ベクトルを生成し、生成した前記周辺ベクトルに対応する前記事例ベクトルを特定するための索引情報を構築する索引構築部と、

任意の特徴ベクトルであるクエリベクトルを指定した検索要求に応じて、前記索引情報を用いて、前記クエリベクトルと完全一致する前記周辺ベクトルに対応する前記事例ベクトルを特定し、特定した前記事例ベクトルに基づく検索結果を出力する検索部と、を備え

前記検索部は、前記索引情報を用いた検索方式と、元の特徴ベクトルまたは縮約ベクトル空間に写像された特徴ベクトルを用いて線形検索する検索方式と、を含む複数の検索方式のうちの任意の検索方式を多段階で組み合わせ、前記事例ベクトルの検索を行う、データ管理システム。

【請求項 10】

前記索引構築部は、前記事例ベクトルを縮約ベクトル空間に写像した縮約事例ベクトルに類似する縮約周辺ベクトルを生成し、生成した前記縮約周辺ベクトルに対応する前記事例ベクトルを特定するための前記索引情報を構築し、

前記検索部は、前記クエリベクトルを前記索引構築部と共通の縮約ベクトル空間に写像した縮約クエリベクトルと完全一致する前記縮約周辺ベクトルに対応する前記事例ベクトルを特定する、請求項 1 乃至 9 のいずれか一項に記載のデータ管理システム。

【請求項 11】

前記事例ベクトルから前記縮約事例ベクトルへの写像、および、前記クエリベクトルから前記縮約クエリベクトルへの写像に、L S H (Locality-Sensitive Hashing) 技術を用いる、請求項 10 に記載のデータ管理システム。

【請求項 12】

前記 L S H 技術として、前記事例ベクトルを各次元が二値のみを取る二値ベクトルである前記縮約事例ベクトルに変換するとともに、前記クエリベクトルを各次元が二値のみを取る二値ベクトルである前記縮約クエリベクトルに変換するビットワイズ L S H を用いる、請求項 11 に記載のデータ管理システム。

【請求項 13】

10

20

30

40

50

前記索引構築部は、二値ベクトルである前記縮約周辺ベクトルを生成し、生成した前記縮約周辺ベクトルを二進数と見立てて整数として格納した前記索引情報を構築し、

前記検索部は、二値ベクトルである前記縮約クエリベクトルを、二進数と見立てて整数として用いて前記縮約クエリベクトルと完全一致する前記縮約周辺ベクトルに対応する前記事例ベクトルを特定する、請求項 1 2 に記載のデータ管理システム。

【請求項 1 4】

前記 L S H 技術として、前記事例ベクトルを整数ベクトルである前記縮約事例ベクトルに変換するとともに、前記クエリベクトルを整数ベクトルである前記縮約クエリベクトルに変換する直積量子化 L S H を用いる、請求項 1 1 に記載のデータ管理システム。

【請求項 1 5】

データ管理システムにおいて実行されるデータ管理方法であって、  
蓄積するデータの特徴ベクトルである事例ベクトルに類似する周辺ベクトルを生成し、生成した前記周辺ベクトルに対応する前記事例ベクトルを特定するための索引情報を構築する索引構築ステップと、

任意の特徴ベクトルであるクエリベクトルを指定した検索要求に応じて、前記索引情報を用いて、前記クエリベクトルと完全一致する前記周辺ベクトルに対応する前記事例ベクトルを特定し、特定した前記事例ベクトルに基づく検索結果を出力する検索ステップと、  
を含み、

前記索引構築ステップでは、少なくとも、前記周辺ベクトルを格納する第 1 列と、該周辺ベクトルに対応する前記事例ベクトルに関する情報を格納する第 2 列とを列要素に持つ

前記検索ステップでは、前記索引を用いて、前記クエリベクトルと完全一致する前記周辺ベクトルに対応する前記テーブルのレコードを求め、求めたレコードの前記第 2 列に格納された情報に基づいて前記事例ベクトルを特定する、データ管理方法。

【請求項 1 6】

コンピュータに、  
蓄積するデータの特徴ベクトルである事例ベクトルに類似する周辺ベクトルを生成し、生成した前記周辺ベクトルに対応する前記事例ベクトルを特定するための索引情報を構築する索引構築部の機能と、

任意の特徴ベクトルであるクエリベクトルを指定した検索要求に応じて、前記索引情報を用いて、前記クエリベクトルと完全一致する前記周辺ベクトルに対応する前記事例ベクトルを特定し、特定した前記事例ベクトルに基づく検索結果を出力する検索部の機能と、  
を実現させるためのプログラムであって、

前記索引構築部は、少なくとも、前記周辺ベクトルを格納する第 1 列と、該周辺ベクトルに対応する前記事例ベクトルに関する情報を格納する第 2 列とを列要素に持つテーブルと、該テーブルにおける前記第 1 列に対する索引とを含む前記索引情報を構築し、

前記検索部は、前記索引を用いて、前記クエリベクトルと完全一致する前記周辺ベクトルに対応する前記テーブルのレコードを求め、求めたレコードの前記第 2 列に格納された情報に基づいて前記事例ベクトルを特定する、プログラム。

【発明の詳細な説明】

【技術分野】

【0001】

本発明の実施形態は、データ管理システム、データ管理方法およびプログラムに関する。

【背景技術】

【0002】

近年、情報通信技術の進展に伴って、多種多様なデータの収集や蓄積が可能となり、ビッグデータ分析やビッグメディア解析などといった大規模データを対象とする情報処理技術が注目を浴びている。こうした大規模データを取り扱うシステムでは、データ規模の加

10

20

30

40

50

速度的な拡大に伴う計算量の肥大化がサービス低下に繋がるため、いかに計算量を削減できるかが重要な課題となっている。

【0003】

データベース検索などのデータ検索では、画像や音楽などのメディア検索を高速に行う方法として、多次元の特徴ベクトルを用いた類似検索が行われる。この類似検索では、特徴ベクトル間の類似度計算を含むベクトル近傍検索、すなわち、ある特徴ベクトル（以下、これを「クエリベクトル」という）に近い特徴ベクトル群を、検索対象となる特徴ベクトル（以下、これを「事例ベクトル」という）群の中から見つけ出す処理が計算量の多くを占めている。このため、ベクトル近傍検索の計算量を削減してデータ検索の実行時間を短縮できるようにすることが求められている。

10

【先行技術文献】

【特許文献】

【0004】

【特許文献1】特開2000-35965号公報

【特許文献2】特開2001-52024号公報

【発明の概要】

【発明が解決しようとする課題】

【0005】

本発明が解決しようとする課題は、ベクトル近傍検索の計算量を削減してデータ検索の実行時間を短縮できるデータ管理システム、データ管理方法およびプログラムを提供することである。

20

【課題を解決するための手段】

【0006】

実施形態のデータ管理システムは、索引構築部と、検索部と、を備える。索引構築部は、蓄積するデータの特徴ベクトルである事例ベクトルに類似する周辺ベクトルを生成し、生成した前記周辺ベクトルに対応する前記事例ベクトルを特定するための索引情報を構築する。検索部は、任意の特徴ベクトルであるクエリベクトルを指定した検索要求に応じて、前記索引情報を用いて、前記クエリベクトルと完全一致する前記周辺ベクトルに対応する前記事例ベクトルを特定し、特定した前記事例ベクトルに基づく検索結果を出力する。前記索引構築部は、少なくとも、前記周辺ベクトルを格納する第1列と、該周辺ベクトルに対応する前記事例ベクトルに関する情報を格納する第2列とを列要素に持つテーブルと、該テーブルにおける前記第1列に対する索引とを含む前記索引情報を構築し、前記検索部は、前記索引を用いて、前記クエリベクトルと完全一致する前記周辺ベクトルに対応する前記テーブルのレコードを求め、求めたレコードの前記第2列に格納された情報に基づいて前記事例ベクトルを特定する。

30

【図面の簡単な説明】

【0007】

【図1】図1は、第1実施形態のデータ管理システムの概要を示すシステム構成図である。

【図2】図2は、データテーブルの具体例を示す図である。

40

【図3】図3は、メディアデータが静止画である場合のデータ登録器による処理手順の一例を示すフローチャートである。

【図4】図4は、メディアデータが動画である場合のデータ登録器による処理手順の一例を示すフローチャートである。

【図5】図5は、索引構築器の構成例を示すブロック図である。

【図6】図6は、LSH即値テーブルの一例を示す図である。

【図7】図7は、LSH即値索引情報生成器による処理手順の一例を示すフローチャートである。

【図8】図8は、LSH近傍展開テーブルの一例を示す図である。

【図9】図9は、LSH近傍展開テーブルを2つに分割することで正規化した例を示す図

50

である。

【図10】図10は、データベース複合索引の一例を示す図である。

【図11】図11は、LSH近傍展開索引情報生成器による処理手順の一例を示すフローチャートである。

【図12】図12は、連想配列と連続メモリ配置型配列を説明する図である。

【図13】図13は、検索器の構成例を示すブロック図である。

【図14】図14は、ベクトル類似性判定部の入出力関係を示す図である。

【図15】図15は、ベクトル類似性判定部による処理手順の一例を示すフローチャートである。

【図16】図16は、厳密検索器による処理手順の一例を示すフローチャートである。

10

【図17】図17は、線形LSH検索器による処理手順の一例を示すフローチャートである。

【図18】図18は、データベース索引LSH検索器による処理手順の一例を示すフローチャートである。

【図19】図19は、データベース索引LSH検索+厳密検索器による処理手順の一例を示すフローチャートである。

【図20】図20は、データベース索引LSH検索+線形LSH検索器による処理手順の一例を示すフローチャートである。

【図21】図21は、第2実施形態の検索器の構成例を示すブロック図である。

【図22】図22は、クエリ摂動型LSH検索器による処理手順の一例を示すフローチャートである。

20

【図23】図23は、第3実施形態の索引構築器の構成例を示すブロック図である。

【図24】図24は、PQLSH近傍展開テーブルの一例を示す図である。

【図25】図25は、PQLSH近傍展開索引情報生成器による処理手順の一例を示すフローチャートである。

【図26】図26は、第3実施形態の検索器の構成例を示すブロック図である。

【図27】図27は、データベース索引PQLSH検索器による処理手順の一例を示すフローチャートである。

【図28】図28は、データ管理システムのハードウェア構成例を示すブロック図である。

30

【発明を実施するための形態】

【0008】

以下、実施形態のデータ管理システム、データ管理方法およびプログラムを、図面を参照して詳細に説明する。

【0009】

実施形態のデータ管理システムは、大規模データを効率よく管理・検索するためのシステムである。データベース管理システムにみられるような大規模データを管理する従来のシステムは、一般的に、下位層としてディスクアクセスなどを最適化するデータ配置機構、上位層として検索条件に基づいて大規模データを高速検索するための索引機構が搭載されている。索引のアルゴリズムには、Bツリーなどの木構造アルゴリズムや、一般的なハッシュアルゴリズムが主に用いられている。

40

【0010】

データベースで利用できる検索条件は、多くの場合、実数・整数・文字列・日付などの基本型（以下、これらを総称して「DB基本型」と呼ぶ）に関する四則演算・集合演算などで構成された論理式である。ベクトルやビット列（すなわち二値ベクトル）同士の類似性については、一部の例外を除き、通常、検索条件として用いることができない。その背景には、ベクトルの類似度計算を含む検索条件を効率化する索引が考案されていないことがある。

【0011】

上述の一部の例外とは、低次元ベクトルを対象とした類似度計算である。多くの著名な

50

データベース管理システムでは、主に時空間データを管理する用途を想定して、2～4次元程度の低次元のベクトルの類似度計算を検索条件に含めることができるようになってい  
る。この類似度計算の実現には、空間木（空間分割法）と呼ばれる索引手法が用いられて  
おり、これにより高速な検索が可能である。ただしベクトルが低次元でない場合は、索引  
サイズの肥大化によって高速化の効果が失われ、実行時間が通常の線形検索と変わらな  
くることが知られている。実施形態で想定する特徴ベクトルは、例えば機械学習用途を想  
定した数百～数億次元の高次元ベクトルであるため、空間木による方法を用いても高速な  
検索は実現できない。

#### 【0012】

しかし、高次元ベクトルの完全一致を高速照合することは可能である。この場合、2つ  
のベクトルの各次元の要素に対する等号条件をANDで結合した条件で検索することと等  
価である。すなわち、この検索条件はスカラ演算のみで構成されており、ベクトルの類似  
度計算は含まれていないため、Bツリーなどの索引を利用できる。

10

#### 【0013】

以上のように、従来一般的なデータベース管理システムは、DB基本型で構成された  
検索条件での大規模データの高速検索が可能であるが、低次元でないベクトル間の類似度  
計算を含むベクトル近傍検索を高速実行することはできない。そこで、以下に示す実施形  
態では、ベクトル近傍検索の計算量を削減することでベクトル近傍検索の高速実行を可能  
とし、データ検索の実行時間を短縮できる新規なデータ管理システムを提案する。

#### 【0014】

なお、ベクトル間の類似度を表す指標としては内積と距離がある。ベクトル間の内積値  
を求める内積計算とベクトル間の距離を求める距離計算は、意味としても計算量としても  
ほぼ同じである。2つのベクトルが類似している場合、内積値は大きくなるが、距離は小  
さくなるという点が異なるだけである。つまり、ベクトル間の内積値が大きいことと、ベ  
クトル間の距離が小さいことは、ともにベクトル間の類似度が高いことを意味する。以下  
では、ベクトル間の類似度を距離で表すものとして説明するが、距離を内積に置き換えて  
もよい。この場合、ベクトル間の距離が小さいほど、つまり、ベクトル間の類似度が高い  
ほど、内積値が大きくなるものと考えればよい。

20

#### 【0015】

ベクトル近傍検索の計算量を削減する方法として、LSH (Locality-Sensitive Hash  
ing) を用いることが考えられる。LSHは、与えられたベクトル群を、離散値のみを取  
る縮約ベクトル空間に写像する技術である（例えば、下記の参考文献1参照）。このLS  
Hによる写像は、写像前の空間におけるベクトル間の距離の相対的な大きさが、写像後の  
空間においてもよく保存されているという性質がある。したがって、写像前のベクトル空  
間でベクトル間距離を計算する代わりに、写像後のベクトル空間でベクトル間距離を計算  
することで、計算を効率化することができる。ただし、距離の大小関係を完全に保存す  
るわけではないため、得られるのは近似解である。

30

参考文献1: Anshumali Shrivastava and Ping Li, "Asymmetric LSH (ALSH)  
for Sublinear Time Maximum Inner Product Search (MIPS)", Advances in Neural Information Processing Systems, 2014.

40

#### 【0016】

LSHアルゴリズムとしては、これまでに、例えば、下記の参考文献2に示されるSimH  
ash (Random projection) や、下記の参考文献3に示されるSpectral Hashingなど、様  
々なアルゴリズムが考案されている。

参考文献2: Moses S. Charikar, "Similarity Estimation Techniques from Rounding Algorithms", Proceedings of the 34<sup>th</sup> Annual ACM Symposium on Theory of Computing, pp. 380-388, doi: 10.1145/509907.509965.

参考文献3: Yair Weiss, Antonio Torralba and Rob Fergus, "Spectral Hashing", Advances in neural information processing systems, 2009.

#### 【0017】

50

L S Hアルゴリズムは、大きく、ビットワイズ方式と直積量子化 (Product Quantization) 方式とに二分できる。ビットワイズ方式のアルゴリズムは、与えられたベクトルに対する写像結果として二値ベクトルを出力する。一方、直積量子化方式のアルゴリズムは、与えられたベクトルに対する写像結果として整数ベクトル (整数の要素からのみ構成される多値ベクトル) を出力する。また、直積量子化方式のアルゴリズムの場合、ベクトルの各次元ごとに、距離に対する重みに相当する情報 (各次元の変化に対してベクトルの距離がどれくらい変化するか) も合わせて生成される。

【 0 0 1 8 】

上述のように、L S Hによって写像されたベクトルを用いることで、ベクトル間距離の計算を効率化することができる。特に二値ベクトルの場合、距離計算処理に、xorなどのビット演算命令、あるいはpopcountなどの専用CPU命令を用いることができるため、大幅な計算量削減が可能である。

10

【 0 0 1 9 】

しかし、L S Hを用いたとしても、ベクトル間距離の計算そのものを回避することはできない。従来のデータベース管理システムでは、上述のように、低次元でないベクトルの距離計算を伴う検索処理に対して有効な索引アルゴリズムが考案されていないため、線形検索を行う必要があり、計算量オーダーをサブリニアにすることができない。その結果、検索対象となる事例ベクトルの数が例えば1億など大規模な場合には、大きな検索時間を要することとなる。そこで、低次元でないベクトルの距離計算を伴う検索処理に対して有効な索引アルゴリズムを導入した新規なデータ管理システムを提案する。

20

【 0 0 2 0 】

ここで、実施形態の基本原理を説明する。実施形態のデータ管理システムは、まず事前処理として、各事例ベクトルのそれぞれについて、事例ベクトルに類似するベクトル、つまり事例ベクトルに対する距離が所定値以下のベクトル (以下、これを「周辺ベクトル」と呼ぶ) 群を生成するとともに、周辺ベクトルから事例ベクトルを特定するための索引情報を構築し、事例ベクトルを含むデータと索引情報とをデータベースに保存する。そして、検索時には、索引情報に基づいて、クエリベクトルと完全一致する周辺ベクトルに対応する事例ベクトルを特定し、特定した事例ベクトルに基づいて、検索要求に対する検索結果を出力する。

【 0 0 2 1 】

事前処理で事例ベクトルに類似する周辺ベクトルを合理的に生成するには、取り扱うベクトルの次元数および各要素値集合のカーディナリティがいずれも小さい必要がある。この条件が満たされなければ、周辺ベクトル数が爆発し、結果として検索時間を短縮できない。この条件を達成するために、上述したL S Hを用いる。L S Hにより生成されるベクトルの要素は、ビットワイズ方式の場合は二値、直積量子化方式の場合は多値を持つ離散値を取り、L S Hの設定にもよるが、多くの場合、要素値集合のカーディナリティおよび次元数が小さい離散ベクトル (縮約ベクトル) が生成される。したがって、基本原理として説明した上述の方法を実行する前に、L S Hを用いて事例ベクトル群およびクエリベクトルを縮約ベクトルに変換しておき、これらを用いて、基本原理として説明した上述の方法による検索を行えばよい。ただし、L S Hを用いるため、検索結果は近似解となる。

30

40

【 0 0 2 2 】

なお、事例ベクトル群やクエリベクトルが元々上記条件を満たすベクトルであるならば、L S Hを用いたベクトル変換を事前に行う必要はなく、基本原理として説明した上述の方法のみを実行すればよい。

【 0 0 2 3 】

また、基本原理として説明した上述の方法では、事例ベクトルに対する距離が所定値以下の周辺ベクトルを生成したが、クエリベクトルに対する距離が所定値以下の周辺ベクトルを生成する構成としてもよい。この場合、検索時に、クエリベクトルに対する距離が所定値以下の周辺ベクトルを生成して、周辺ベクトルと完全一致する事例ベクトルを特定し、特定した事例ベクトルに基づいて、検索要求に対する検索結果を出力する。

50

## 【 0 0 2 4 】

以下では、ベクトル間距離の計算が検索条件に含まれるクエリを処理してメディア検索を行うデータ管理システムへの適用例について具体的に説明する。メディア検索の応用題材としては、顔画像（映像含む）検索を取り上げる。顔画像検索は、システムに予め登録された画像・映像などのメディア群の中から、クエリとして与えられた画像・映像に含まれる顔と似た顔を含む画像・映像の箇所を見つけ出す処理である。なお、ここでは応用題材として顔画像検索を取り上げるが、実施形態のデータ管理システムは、例えば、音楽検索、物体画像検索、シーン画像検索、テキスト意味検索、センサパタン検索、株価パタン検索、電力使用パタン検索など、様々なメディア検索やセンサデータ検索に応用することができる。

10

## 【 0 0 2 5 】

## &lt; 第 1 実施形態 &gt;

まず、第 1 実施形態のデータ管理システムについて説明する。本実施形態では、写像処理にはビットワイズ L S H を用いるものとし、写像後の二値ベクトル（ハッシュ）間の距離としてハミング距離を用いるものとする。ただし、上述したように、取り扱う事例ベクトルおよびクエリベクトルが元々低次元の二値ベクトルであるならば、ビットワイズ L S H による写像は不要であり、後述の事例ハッシュ（縮約事例ベクトル）として事例ベクトルをそのまま用い、クエリハッシュ（縮約クエリベクトル）としてクエリベクトルをそのまま用いればよい。また、この場合、後述の周辺ハッシュ（縮約周辺ベクトル）として、事例ベクトルを元に生成される周辺ベクトルを用いればよい。

20

## 【 0 0 2 6 】

図 1 は、第 1 実施形態のデータ管理システムの概要を示すシステム構成図である。本実施形態のデータ管理システムは、図 1 に示すように、データ登録器 1 0 0 と、索引構築器 2 0 0 と、検索器 3 0 0 と、データベース 4 0 0 とを備える。

## 【 0 0 2 7 】

データ登録器 1 0 0 は、顔を含む画像や映像などのメディアデータ群 1 0 を外部から受け取り、受け取ったメディアデータ群 1 0 を対象として顔に関する解析を行い、顔特徴に関する特徴ベクトル（事例ベクトル）群を含むデータテーブル 2 0 をデータベース 4 0 0 内に生成する。

## 【 0 0 2 8 】

索引構築器 2 0 0 は、索引構築命令文 3 0 を外部から受け取ると、受け取った索引構築命令文 3 0 に対応するデータテーブル 2 0 をデータベース 4 0 0 から取り出し、索引構築命令文 3 0 の命令に従って、顔特徴ベクトルの索引として用いる索引情報 4 0 をデータベース 4 0 0 内に生成する。

30

## 【 0 0 2 9 】

検索器 3 0 0 は、顔特徴に関するクエリベクトルを含む拡張 S Q L 5 0（検索要求の一例）を外部から受け取ると、データベース 4 0 0 から該当するデータテーブル 2 0 および索引情報 4 0 を取得し、拡張 S Q L 5 0 に含まれるクエリベクトルと類似する顔特徴ベクトルが含まれる検索結果データセット 6 0 を出力する。

## 【 0 0 3 0 】

まず、データ登録器 1 0 0 の詳細を説明する。図 2 は、データ登録器 1 0 0 がデータベース 4 0 0 内に生成するデータテーブル 2 0 の具体例を示す図である。データ登録器 1 0 0 が外部から受け取るメディアデータ群 1 0 に含まれるメディアデータは、dir/a.jpg、dir/b.jpg、dir/a.mpg の 3 種類のファイルであるものとする。このうち、dir/a.jpg、dir/b.jpg は静止画ファイルであり、dir/a.mpg は動画ファイルである。

40

## 【 0 0 3 1 】

データ登録器 1 0 0 は、外部からメディアデータ群 1 0 を受け取ると、このメディアデータ群 1 0 に含まれるメディアデータそれぞれについて、以下に示す処理を実行する。

## 【 0 0 3 2 】

図 3 は、メディアデータが静止画である場合のデータ登録器 1 0 0 による処理手順の一

50

例を示すフローチャートである。データ登録器100は、メディアデータが静止画である場合、例えば以下のステップS101～ステップS105の処理を実行して、メディアデータをデータテーブル20に登録する。

【0033】

ステップS101：データ登録器100は、入力された静止画に対して顔領域抽出処理を行う。この結果、入力された静止画から、顔が写っている可能性の高い画像領域（顔画像領域）がすべて抽出される。なお、顔領域抽出処理は既存技術を用いればよいため、ここでは詳細な説明を省略する。

【0034】

ステップS102：データ登録器100は、ステップS101で抽出した顔画像領域を順に取り出す。

10

【0035】

ステップS103：データ登録器100は、ステップS102で取り出した顔画像領域に対して特徴生成処理を行って特徴ベクトル（事例ベクトル）を得る。なお、特徴生成処理には既存技術を用いればよいため、ここでは詳細な説明を省略する。

【0036】

ステップS104：データ登録器100は、ステップS103で得られた特徴ベクトル、入力された静止画に与えられたメディアデータ名、顔画像領域の座標の3つの情報を含むレコードをデータテーブル20に追加する。

【0037】

20

ステップS105：データ登録器100は、ステップS101で抽出した顔画像領域をすべて取り出したか否かを判定する。そして、判定の結果がNoであればステップS102に戻って以降の処理を繰り返し、判定の結果がYesであれば処理を終了する。

【0038】

図4は、メディアデータが動画である場合のデータ登録器100による処理手順の一例を示すフローチャートである。データ登録器100は、メディアデータが動画である場合、例えば以下のステップS201～ステップS212の処理を実行して、メディアデータをデータテーブル20に登録する。

【0039】

ステップS201：データ登録器100は、入力された動画から先頭フレーム画像を取り出す。

30

【0040】

ステップS202：データ登録器100は、ステップS201で取り出した先頭フレーム画像に対して顔領域抽出処理および特徴生成処理を行う。この結果、顔が写っている可能性の高い画像領域（顔画像領域）が確信度付ですべて抽出され、それぞれの顔画像領域から特徴ベクトルが得られる。

【0041】

ステップS203：データ登録器100は、入力された動画から次のフレーム画像を取り出す。

【0042】

40

ステップS204：データ登録器100は、ステップS203で取り出したフレーム画像に対して、ステップS202と同様の顔領域抽出処理および特徴生成処理を行う。

【0043】

ステップS205：データ登録器100は、当該フレームについて、前フレームとの間で顔追跡処理を行う。顔追跡処理とは、フレーム間で画像と領域座標が類似する顔画像領域のペアを見つけ出し、これらを同一被写体と解釈する処理である。なお、顔追跡処理は既存技術を用いればよいため、ここでは詳細な説明を省略する。

【0044】

ステップS206：データ登録器100は、ステップS205の顔追跡処理の結果、同一被写体と判定された顔画像領域を同一グループとして束ねる。

50

## 【 0 0 4 5 】

ステップ S 2 0 7 : データ登録器 1 0 0 は、次のフレームがあるか否かを判定する。そして、判定の結果が Y e s であればステップ S 2 0 3 に戻って以降の処理を繰り返し、N o であればステップ S 2 0 8 に進む。

## 【 0 0 4 6 】

ステップ S 2 0 8 : データ登録器 1 0 0 は、以上の処理により生成されたグループを順に取り出す。

## 【 0 0 4 7 】

ステップ S 2 0 9 : データ登録器 1 0 0 は、ステップ S 2 0 8 で取り出したグループに含まれるフレームのうち、最も早く出現したフレームの出現時刻である第 1 出現時刻と、最も遅く出現したフレームの出現時刻である第 2 出現時刻とを取得する。

10

## 【 0 0 4 8 】

ステップ S 2 1 0 : データ登録器 1 0 0 は、ステップ S 2 0 8 で取り出したグループに含まれるフレームのうち最も確信度の高いフレームの出現時刻である第 3 出現時刻と、特徴ベクトルと、顔画像領域の座標とを取得する。

## 【 0 0 4 9 】

ステップ S 2 1 1 : データ登録器 1 0 0 は、ステップ S 2 0 9 で取得した第 1 出現時刻および第 2 出現時刻と、ステップ S 2 1 0 で取得した第 3 出現時刻、特徴ベクトル、および顔画像領域の座標と、入力された動画に与えられたメディアデータ名との 6 つの情報を含むレコードをデータテーブル 2 0 に追加する。

20

## 【 0 0 5 0 】

ステップ S 2 1 2 : データ登録器 1 0 0 は、生成されたグループをすべて取り出したか否かを判定する。そして、判定の結果が N o であればステップ S 2 0 8 に戻って以降の処理を繰り返し、判定の結果が Y e s であれば処理を終了する。

## 【 0 0 5 1 】

次に、索引構築器 2 0 0 の詳細を説明する。索引構築器 2 0 0 がデータベース 4 0 0 内に生成する索引情報 4 0 は、特徴ベクトルを列に持つデータテーブル 2 0 に対するベクトル近傍検索を高速化するための補助情報である。索引構築器 2 0 0 は、入力された索引構築命令文 3 0 のタイプに応じて、L S H 即値索引情報と L S H 近傍展開索引情報との 2 種類の索引情報 4 0 を生成することができる。この索引情報 4 0 には、後述するテーブルデータとデータベース索引 ( B ツリーなど ) が含まれる。

30

## 【 0 0 5 2 】

索引構築器 2 0 0 は、データテーブル 2 0 の特徴ベクトルを格納する同一の行に対してハッシュのビット長の指定が異なる複数の索引構築命令文 3 0 を実行することにより、複数の索引情報 4 0 を構築することができる。例えば、ビット長が短いハッシュを用いた索引情報 4 0 とビット長が長いハッシュを用いた索引情報 4 0 との 2 種類を構築しておけば、精度より速度を重視した検索を行う場合は前者の索引情報 4 0 を用い、速度より精度を重視した検索を行う場合は後者の索引情報 4 0 を用いるといったように、検索時に用途に応じた使い分けができるようになる。

## 【 0 0 5 3 】

図 5 は、索引構築器 2 0 0 の構成例を示すブロック図である。索引構築器 2 0 0 は、図 5 に示すように、索引構築命令タイプ分類器 2 1 0 と、L S H 即値索引情報生成器 2 2 0 と、L S H 近傍展開索引情報生成器 2 3 0 とを備える。

40

## 【 0 0 5 4 】

索引構築命令タイプ分類器 2 1 0 は、索引構築器 2 0 0 に入力された索引構築命令文 3 0 を、そのタイプに応じて、L S H 即値索引情報作成器 2 2 0 と L S H 近傍展開索引情報生成器 2 3 0 とのいずれかに受け渡すスイッチャである。

## 【 0 0 5 5 】

索引構築器 2 0 0 に入力される索引構築命令文 3 0 は、例えば、以下のような構成の命令文である。ただし xxxx 部には何らかの値が入る。

50

Table:xxxx,Column:xxxx,Algo:xxxx,BitLen:xxxx,HammingDist:xxxx

【 0 0 5 6 】

上記構成の索引構築命令文 3 0 は、Table項目で、索引情報 4 0 の対象となるデータテーブル 2 0 のテーブル名が指定され、Column項目で、索引情報 4 0 の対象となるデータテーブル 2 0 内の列名が指定される。この列名が示すデータテーブル 2 0 内の列には、特徴ベクトルが格納されている必要がある。

【 0 0 5 7 】

また、Algo項目では、索引情報 4 0 の生成に用いる L S H アルゴリズムのアルゴリズム名が指定される。アルゴリズム名は、システム内に実装されている L S H アルゴリズムの中から選ばれる。例えば、システム内にSimHashという名前の L S H アルゴリズムとSpectralHashという名前の L S H アルゴリズムのみが実装されているとすれば、これらのいずれかのアルゴリズム名が指定される。本実施形態では、上述したように、ビットワイズ L S H を用いるため、ビットワイズ L S H のアルゴリズム名が指定される。

【 0 0 5 8 】

Bitlen項目では、Algo項目で選んだ L S H アルゴリズムが出力するハッシュのビット長が指定される。HammingDist項目では、L S H 近傍展開索引情報生成器 2 3 0 において扱われるハミング距離の上限（後述の周辺ハッシュの範囲）が指定される。例えばHammingDist項目で“ 2 ”が指定されている場合、L S H 近傍展開索引情報生成器 2 3 0 において扱われるハミング距離は“ 0 ”と“ 1 ”と“ 2 ”である。なお、HammingDist項目は省略してもよい。

【 0 0 5 9 】

索引構築命令タイプ分類器 2 1 0 は、以上のように構成される索引構築命令文 3 0 が索引構築器 2 0 0 に入力されると、この索引構築命令文 3 0 を、L S H 即値索引情報生成器 2 2 0 と L S H 近傍展開索引情報生成器 2 3 0 のいずれかに渡す。例えば、索引構築命令タイプ分類器 2 1 0 は、索引構築器 2 0 0 に入力された索引構築命令文 3 0 にHammingDist項目がなければ、この索引構築命令文 3 0 を L S H 即値索引情報生成器 2 2 0 に渡し、HammingDist項目があれば、この索引構築命令文 3 0 を L S H 近傍展開索引情報生成器 2 3 0 に渡す。

【 0 0 6 0 】

L S H 即値索引情報生成器 2 2 0 は、索引構築命令タイプ分類器 2 1 0 から索引構築命令文 3 0 を受け取った場合に、この索引構築命令文 3 0 に従って、L S H 即値テーブルと、この L S H 即値テーブルのHashValue列に対するデータベース索引とを含む L S H 即値索引情報 4 0 A を生成する。

【 0 0 6 1 】

図 6 は、L S H 即値テーブルの一例を示す図である。L S H 即値テーブルは、図 6 に示すように、HashValue列（型：blob）のみからなるテーブルデータである。L S H 即値索引情報生成器 2 2 0 は、索引構築命令文 3 0 で指定されたデータテーブル 2 0 の各レコードに対して、指定された列に格納される特徴ベクトルを順に取り出し、指定された L S H アルゴリズムを用いて、指定された出力ビット長のハッシュを生成し、得られた各ハッシュのみを持つレコードを L S H 即値テーブルに登録する。

【 0 0 6 2 】

なお、ハッシュは本質的には二値ベクトルであるが、本実施形態ではハッシュをブール値の配列として保存するのではなく、二進数と見立てて整数として保存する。すなわち、L S H 即値テーブルのHashValue列は、整数または整数配列（計算機アーキテクチャのビット数上限を超えた場合の処置）を格納するblobである。これにより、保存領域サイズを削減することができる。

【 0 0 6 3 】

なお、この L S H 即値テーブルをデータテーブル 2 0 に組み込む構成としてもよい。このような構成とすることにより、検索全体の処理においてサブクエリ処理を 1 回減らすことができる。L S H 即値テーブルをデータテーブル 2 0 に組み込むためには、上述のデー

10

20

30

40

50

タ登録器100がデータテーブル20を生成する際に、このデータテーブル20にハッシュを格納する列を予め用意しておく必要がある。また、上述したように、データテーブル20の特徴ベクトルを格納する同一の行に対して複数の索引情報40を構築できるようにするためには、データテーブル20内にビット長が異なる複数種類のハッシュを格納する複数の列を確保しておく必要がある。このような構成とした場合、LSH即値索引情報生成器220は、特徴ベクトルを元に生成したハッシュをLSH即値テーブルに登録する代わりに、データテーブル20の該当する列に格納すればよい。

【0064】

LSH即値索引情報生成器220は、以上のようなLSH即値テーブルの生成と合せて、LSH即値テーブルのHashValue列に対するデータベース索引を構築し、これらをLSH即値索引情報40Aとしてデータベース400内に保存する。

10

【0065】

図7は、LSH即値索引情報生成器220による処理手順の一例を示すフローチャートである。LSH即値索引情報生成器220は、索引構築命令タイプ分類器210から索引構築命令文30を受け取ると、例えば以下のステップS301～ステップS307の処理を実行し、LSH即値テーブルとデータベース索引とを含むLSH即値索引情報40Aをデータベース400に保存する。

【0066】

ステップS301：LSH即値索引情報生成器220は、データベース400内に、名前：HashValue・型：blobの列のみを持つテーブル（空のLSH即値テーブル）を生成する。なお、テーブル名の形式は「AAAA\_BBBB\_CCCC\_DDDD」とする。ただし、AAAAは処理対象のデータテーブル20のテーブル名、BBBBは処理対象の列の列名、CCCCはLSHアルゴリズムのアルゴリズム名、DDDDはビット長とする。例えば、aTable\_feat\_SimHash\_64といったテーブル名が生成したテーブルに与えられる。

20

【0067】

ステップS302：LSH即値索引情報生成器220は、索引構築命令文30で指定されたデータテーブル20から、索引構築命令文30で指定された列の特徴ベクトルを順に取り出す。

【0068】

ステップS303：LSH即値索引情報生成器220は、ステップS302で取り出した特徴ベクトルに対して、索引構築命令文30で指定されたLSHアルゴリズムを用いて、索引構築命令文30で指定された出力ビット長のハッシュを生成する。

30

【0069】

ステップS304：LSH即値索引情報生成器220は、HashValue列の値としてステップS303で得られた事例ハッシュを持つレコードを、LSH即値テーブルに追加する。

【0070】

ステップS305：LSH即値索引情報生成器220は、処理対象のデータテーブル20から特徴ベクトルをすべて取り出したか否かを判定する。そして、判定の結果がNoであればステップS302に戻って以降の処理を繰り返し、YesであればステップS306に進む。

40

【0071】

ステップS306：LSH即値索引情報生成器220は、以上の処理を経て生成されたLSH即値テーブルのHashValue列に対するデータベース索引（Bツリーなど）を構築し、データベース400に保存する。

【0072】

ステップS307：LSH即値索引情報生成器220は、ハッシュの生成に用いたLSHアルゴリズムのモデルパラメタセットなど（例えばSimHashにおける射影ベクトル群など）があれば、テーブル名と同じファイル名（例えばaTable\_feat\_SimHash\_64）でデータベース400に保存する。

50

## 【 0 0 7 3 】

L S H近傍展開索引情報生成器 2 3 0 は、索引構築命令タイプ分類器 2 1 0 から索引構築命令文 3 0 を受け取った場合に、この索引構築命令文 3 0 に従って、L S H近傍展開テーブルと、このL S H近傍展開テーブルのHashValue列およびHammingDistance列に対するデータベース複合索引（コンポジットインデックス）とを含むL S H近傍展開索引情報 4 0 Bを生成する。

## 【 0 0 7 4 】

図 8 は、L S H近傍展開テーブルの一例を示す図である。L S H近傍展開テーブルは、図 8 に示すように、HashValue列（型：blob）、HammingDistance列（型：integer）およびDataRowID列（型：integer）からなるテーブルデータである。L S H近傍展開索引情報生成器 2 3 0 は、索引構築命令文 3 0 で指定されたデータテーブル 2 0 から特徴ベクトルを取り出し、各特徴ベクトル（事例ベクトル）について、索引構築命令文 3 0 で指定されたL S Hアルゴリズムおよびビット長に従ってハッシュ（縮約事例ベクトル）を生成する。これを「事例ハッシュ」と呼ぶ。次に、L S H近傍展開索引情報生成器 2 3 0 は、生成した各事例ハッシュについて、索引構築命令文 3 0 で指定されたハミング距離以内のハッシュ（縮約周辺ベクトル）群を生成する。ここで生成されるハッシュを「周辺ハッシュ」と呼ぶ。そして、L S H近傍展開索引情報生成器 2 3 0 は、例えば、周辺ハッシュと、周辺ハッシュの事例ハッシュからのハミング距離と、事例ハッシュの元になる事例ベクトルが属するデータテーブル 2 0 のレコードの行 ID との 3 つ組からなるレコードを、L S H近傍展開テーブルに登録する。

## 【 0 0 7 5 】

なお、本実施形態では、後述のデータベース索引L S H検索処理において、検索条件として用いるハミング距離を 0 から順次インクリメントして検索結果を追加する処理を行う。このため、L S H近傍展開テーブルの第 2 列にHammingDistance列を設け、周辺ハッシュの事例ハッシュからのハミング距離をこのHammingDistance列に格納している。しかし、検索時にL S H近傍展開テーブルの距離上限を検索条件として用いて検索を行う構成とすることも可能である。この場合、L S H近傍展開テーブルにHammingDistance列を設ける必要はない。このような構成とした場合は、後述のデータベース索引L S H検索処理において、出力する検索結果データセットの件数をハミング距離を利用して制御することはできなくなるが、L S H近傍展開テーブルのサイズを小さくできる効果がある。また、このような構成とした場合は、L S H近傍展開テーブルのHashValue列およびHammingDistance列に対するデータベース複合索引を構築する代わりに、L S H近傍展開テーブルのHashValue列に対するデータベース索引を構築すればよい。

## 【 0 0 7 6 】

また、図 8 に示したL S H近傍展開テーブルには、HashValue列に格納される周辺ハッシュの値として重複した値が多く含まれているため、正規化を行うことでコンパクト化することができる。

## 【 0 0 7 7 】

図 9 は、図 8 に示したL S H近傍展開テーブルを 2 つに分割することで正規化した例を示す図である。L S H近傍展開テーブルを図 9 のような構成とした場合、後述のデータベース索引L S H検索処理において、HashValueとHammingDistanceに関する検索条件を用いて図 9 ( a ) のテーブルの行 ID を確定し、当該行 ID を元に生成した図 9 ( b ) のテーブルのTable1\_RowIDに関する検索条件を用いてDataRowIDを見つければよい。なお、図 9 ( b ) のテーブルの第 1 列であるTable1\_RowIDが図 9 ( a ) のテーブルの対応する行の行 ID を格納しており、このリンクを用いて図 9 ( a ) のテーブルと図 9 ( b ) のテーブルとを結合（JOIN）すれば、図 8 に示したL S H近傍展開テーブルになる。

## 【 0 0 7 8 】

また、図 9 ( a ) のテーブルと図 9 ( b ) のテーブルのそれぞれを、連想配列を用いて実装してもよい。連想配列として実装する場合、図 9 ( a ) の連想配列のキーはHashValueとHammingDistanceを組み合わせたものであり、値は 1 つ以上の行 ID である。また、図

9 ( b ) の連想配列のキーはTable1\_RowIDであり、値は1つ以上のDataRowIDである。

【 0 0 7 9 】

図10は、LSH近傍展開テーブルのHashValue列およびHammingDistance列に対するデータベース複合索引の一例を示す図である。このデータベース複合索引は、図10に示すように、HashValueおよびHammingDistanceの値に沿って分岐する木構造となっており、葉に該当するLSH近傍展開テーブルの行IDが記載されている。このデータベース複合索引を用いることで、検索時間を $O(N)$ から $O(\log(N))$ に短縮させることができる。

【 0 0 8 0 】

なお、データベース複合索引の構成を少し変更することで、LSH近傍展開テーブルを不要とすることもできる。具体的には、データベース複合索引の葉に、LSH近傍展開テーブルの行IDを格納する代わりにDataRowIDの値を格納する構成とする。これにより、DataRowIDしか検索できなくなるが、LSH近傍展開テーブルをデータベース400に保存しておく必要がなくなる。この場合、LSH近傍展開索引情報40Bはデータベース複合索引のみとなる。

10

【 0 0 8 1 】

図11は、LSH近傍展開索引情報生成器230による処理手順の一例を示すフローチャートである。LSH近傍展開索引情報生成器230は、索引構築命令タイプ分類器210から索引構築命令文30を受け取ると、例えば以下のステップS401～ステップS410の処理を実行し、LSH近傍展開テーブルとデータベース複合索引とを含むLSH近傍展開索引情報40Bをデータベース400に保存する。

20

【 0 0 8 2 】

ステップS401：LSH近傍展開索引情報生成器230は、データベース400内に、名前：HashValue・型：blobの列と、名前：HammingDistance・型：integerの列と、名前：DataRowID・型：integerの列とを持つテーブル（空のLSH近傍展開テーブル）を生成する。なお、テーブル名の形式は「AAAA\_BBBB\_CCCC\_DDDD」とする。ただし、AAAAは処理対象のデータテーブル20のテーブル名、BBBBは処理対象の列の列名、CCCCはLSHアルゴリズムのアルゴリズム名、DDDDはビット長とする。例えば、aTable\_feat\_PQ\_3といったテーブル名が生成したテーブルに与えられる。

30

【 0 0 8 3 】

ステップS402：LSH近傍展開索引情報生成器230は、索引構築命令文30で指定されたデータテーブル20から、索引構築命令文30で指定された列の特徴ベクトルを順に取り出す。

【 0 0 8 4 】

ステップS403：LSH近傍展開索引情報生成器230は、ステップS402で取り出した特徴ベクトル（事例ベクトル）に対して、索引構築命令文30で指定されたLSHアルゴリズムを用いて、索引構築命令文30で指定されたビット長の事例ハッシュを生成する。

【 0 0 8 5 】

ステップS404：LSH近傍展開索引情報生成器230は、ステップS403で得られた事例ハッシュからの距離が指定距離以下の周辺ハッシュをすべて生成する。

40

【 0 0 8 6 】

ステップS405：LSH近傍展開索引情報生成器230は、ステップS404で得られた周辺ハッシュを順に取り出す。

【 0 0 8 7 】

ステップS406：LSH近傍展開索引情報生成器230は、ステップS405で取り出した周辺ハッシュ、この周辺ハッシュの事例ハッシュからのハミング距離、事例ハッシュの元になる事例ベクトルが属するデータテーブル20内のレコードの行IDの3つ組からなるレコードを、LSH近傍展開テーブルに追加する。

【 0 0 8 8 】

50

ステップS 4 0 7 : L S H近傍展開索引情報生成器 2 3 0 は、ステップS 4 0 4 で得られた周辺ハッシュをすべて取り出したか否かを判定する。そして、判定の結果がN oであればステップS 4 0 5 に戻って以降の処理を繰り返し、Y e sであればステップS 4 0 8 に進む。

【 0 0 8 9 】

ステップS 4 0 8 : L S H近傍展開索引情報生成器 2 3 0 は、処理対象のデータテーブル 2 0 から特徴ベクトルをすべて取り出したか否かを判定する。そして、判定の結果がN oであればステップS 4 0 2 に戻って以降の処理を繰り返し、Y e sであればステップS 4 0 9 に進む。

【 0 0 9 0 】

ステップS 4 0 9 : L S H近傍展開索引情報生成器 2 3 0 は、以上の処理を経て生成されたL S H近傍展開テーブルのHashValue列およびHammingDistance列に対するデータベース複合索引(コンポジットインデックス)を構築し、データベース 4 0 0 に保存する。

【 0 0 9 1 】

ステップS 4 1 0 : L S H近傍展開索引情報生成器 2 3 0 は、事例ハッシュの生成に用いたL S Hアルゴリズムのモデルパラメタセットなどがあれば、テーブル名と同じファイル名(例えばaTable\_feat\_PQ\_3)でデータベース 4 0 0 に保存する。

【 0 0 9 2 】

ここで、上記ステップS 4 0 4 における周辺ハッシュの生成方法の具体例について説明する。以下では、周辺ハッシュの起点となる事例ハッシュのビット列が“ 0 1 0 ”であるものとし、ハミング距離 2 以内の周辺ハッシュを列挙するという設定例で説明する。

【 0 0 9 3 】

まず、ハミング距離 0 の周辺ハッシュは、事例ハッシュと同じハッシュである。すなわち、事例ハッシュのビット列が“ 0 1 0 ”であれば、この事例ハッシュに対してハミング距離 0 の周辺ハッシュは“ 0 1 0 ”のみである。

【 0 0 9 4 】

次に、ハミング距離 1 の周辺ハッシュは、事例ハッシュの任意のビットを 1 つだけ反転させたものである。事例ハッシュのビット長が 3 であれば、反転させるビットの選択方法は、第 1 ビットのみ、第 2 ビットのみ、第 3 ビットをみの 3 種類である。事例ハッシュのビット列が“ 0 1 0 ”であれば、各選択方法でビット反転を行った結果は“ 1 1 0 ”、“ 0 0 0 ”および“ 0 1 1 ”である。これら 3 つのハッシュが、事例ハッシュ“ 0 1 0 ”に対してハミング距離が 1 の周辺ハッシュである。

【 0 0 9 5 】

次に、ハミング距離 2 の周辺ハッシュは、事例ハッシュの任意のビットを 2 つ反転させたものである。事例ハッシュのビット長が 3 であれば、反転させるビットの選択方法は、第 1 ビット+第 2 ビット、第 1 ビット+第 3 ビット、第 2 ビット+第 3 ビットの 3 種類である。事例ハッシュのビット列が“ 0 1 0 ”であれば、各選択方法でビット反転を行った結果は“ 1 0 0 ”、“ 1 1 1 ”および“ 0 0 1 ”である。これら 3 つのハッシュが、事例ハッシュ“ 0 1 0 ”に対してハミング距離が 2 の周辺ハッシュである。

【 0 0 9 6 】

以上をまとめると、事例ハッシュ“ 0 1 0 ”に対するハミング距離 0 の周辺ハッシュは“ 0 1 0 ”、ハミング距離 1 の周辺ハッシュは“ 1 1 0 ”、“ 0 0 0 ”および“ 0 1 1 ”、ハミング距離 2 の周辺ハッシュは“ 1 0 0 ”、“ 1 1 1 ”および“ 0 0 1 ”である。このように、低コストな処理で周辺ハッシュをすべて生成することができる。

【 0 0 9 7 】

なお、以上の方法により事例ハッシュに対する周辺ハッシュを生成する際、ビットの重要度などを元に、反転を許可するビットを制限する仕組みを備えてもよい。例えば、上記参考文献 3 のように主成分分析をベースとしたハッシュ生成法の場合、各ビットに対応する固有値から当該ビットの弁別性(重要度)を割り出すことができる。事例ハッシュにおける弁別性の高いビットを反転して生成した周辺ハッシュは、弁別性が失われた周辺ハッ

10

20

30

40

50

シュとなるため、このような周辺ハッシュをLSH近傍展開テーブルに追加してもヒットする確率が低いと考えられる。したがって、このような弁別性の高いビットは反転させないように制限を加えるようにしてもよい。このように、反転を許可するビットを制限することにより、LSH近傍展開テーブルの記憶領域サイズを小さくすることができる。

#### 【0098】

LSH近傍展開テーブルは、検索キーとなる列と検索結果となる列が決まっているため、連想配列（キーバリューストア、C++標準ライブラリであればstd::mapやstd::unordered\_mapなど）で簡潔に実装することができる。

#### 【0099】

すなわち、検索結果はDataRowIDであるため、連想配列の値としてDataRowIDの可変長の連続メモリ配置型配列（C++標準ライブラリであればstd::vectorなど）を用いる。検索条件はHashValueとHammingDistanceのANDであるため、これら2つの値を用いてユニークになる連想配列のキーを生成する。図8に例示したLSH近傍展開テーブルの場合、HashValue列は二進数表記で“000”～“111”の値を取り、HammingDistanceは“0”～“1”の値を取るため、キー値を以下の算出式で決定するなどが考えられる。

キー値 = HashValue列の値 × 2 + HammingDistance列の値

#### 【0100】

図12は、連想配列と連続メモリ配置型配列を説明する図である。連想配列は、キーから値へ対応付ける方法として、図12(a)に示すように、データベース索引のように木やハッシュなどのデータ構造を内部に持つ。上記のキーが定義域においてほぼ充当されているならば、連想配列の代わりに、図12(b)に示すような連続メモリ配置型配列（C++標準ライブラリであればstd::vectorなど）を用いてもよい。連続メモリ配置型配列として確保する要素数は、HashValueとHammingDistanceの全組合せ数とする。連続メモリ配置型配列を用いたデータ構造は、連想配列を用いたデータ構造と比べ、キーの充当率が高い場合に限りメモリ使用量が低減され、さらに処理速度が高速化される。

#### 【0101】

次に、検索器300の詳細を説明する。この検索器300の機能概要および特徴は以下の通りである。この検索器300では、拡張SQL50を用いて、ベクトル類似性に関する0個以上の条件を含む論理条件でのレコード検索ができる。これにより、データベース利用者は、一般的なデータ照合と特徴ベクトル類似性判定とを条件として結び付けた論理式での検索が行えるようになる。

#### 【0102】

また、ベクトル類似性判定（ベクトル近傍検索）の処理は、検索対象のレコード数に応じて、5種類の方式を内部的に切り替えることができる。この切り替えは事前に設定されたシステム設定に基づいて行われる。5種類の方式には速度、精度、件数制御の有無について長短がある。これら5種類の方式の1つであるデータベース索引LSH検索は、クエリに対して、検索時にベクトル距離計算を行わずデータベース索引のみを用いてベクトル類似性判定（ベクトル近傍検索）を実現する方式である。この方式は、検索時間を大幅に削減する効果を持つ。

#### 【0103】

図13は、検索器300の構成例を示すブロック図である。検索器300は、拡張SQL50を入力として受け取り、拡張SQL50に従った検索処理を実行して検索結果データセット60を出力するものであり、図13に示すように、検索条件処理部310と、検索出力部320とを備える。検索条件処理部310は、内部に様々な検索条件に対する処理を行う機能モジュールを持つ。このうち、特に本実施形態に特徴的なベクトル類似性に関する検索条件に対する処理を行う機能モジュールがベクトル類似性判定処理部330である。また、それ以外の検索条件に対する処理を行う機能モジュールを、他の検索条件処理部340と総称する。検索条件処理部310は、これらベクトル類似性判定処理部330の検索結果と他の検索条件処理部340との集合演算を行い、検索条件全体の処理を完成させる。検索出力部320は、検索条件処理部310の処理の結果に基づいてデータベ

10

20

30

40

50

ース400からレコードを取り出し、検索結果データセット60として出力する。

【0104】

なお、検索器300による検索処理の大半は従来のSQLの検索処理と同様であり、検索条件処理部310がSQLにおけるfrom命令部（データソースの特定処理）およびwhere命令部（検索条件処理）にあたり、検索出力部320がSQLにおけるselect命令部にあたる。本実施形態における固有の処理は、検索条件処理部310の中の1つであるベクトル類似性判定処理部330が実行する処理である。このため、以下ではベクトル類似性判定処理部330を中心に説明し、従来技術をそのまま適用できる部分については適宜説明を省略する。

【0105】

まず、検索器300が入力として受け取る拡張SQL50について説明する。拡張SQL50は、従来のSQLに、ベクトル類似性に関する条件にあたるvnn関数を記述できるよう拡張した問合せ言語である。vnn関数は、以下の例のように、where命令内で一般的なDB基本型データへの条件群と組合せ可能な項として用いる。

```
select * from aTable where vnn(featsimhash_64,(10,20,30),10) and annualincome > 10000000
```

【0106】

vnn関数の仕様は以下の通りである。

入力：vnn関数が検索対象とするデータセットを入力とする。この入力は、検索条件処理部310が内部的に生成してvnn関数を呼び出す際に渡す。通常の構文であれば、from命令で指定されたテーブルの全レコード群である。

第1引数：第1引数は、事例ベクトル（検索対象とする特徴ベクトル）を格納する列名と、その列に適用されているLSHアルゴリズムのアルゴリズム名と、その列に適用されているLSHアルゴリズムの出力ビット長（後述の第3実施形態ではより一般化してオプションパラメタとしている）とをアンダースコアで結んだテキストである。

第2引数：第2引数は、クエリベクトルである。

第3引数：第3引数は、出力上位件数（出力対象とするベクトルの上位件数）である。ただし、クエリベクトルとの距離の小さいベクトルを上位とみなすものとする。

出力：検索条件処理部310（from, whereなどの処理）に対して、引数で示された条件に合致する検索結果を返す。

なお、ここでは出力条件にあたる第3引数として出力上位件数を与えるものとしているが、これに代えて距離上限を第3引数として与えるようにしてもよい。

【0107】

図14は、ベクトル類似性判定部330の入出力関係を示す図である。ベクトル類似性判定部330は、図14に示すように、検索対象データセット70を入力として受け取り、vnn関数で表現された検索条件に基づくベクトル類似性判定（ベクトル近傍検索）処理を行って、検索結果80を返す。本実施形態では、ベクトル類似性判定部330による検索方式として、厳密検索（Strict Search）と、線形LSH検索（Linear LSH Search）と、データベース索引LSH検索（DB-Indexed LSH Search）と、データベース索引LSH検索+厳密検索と、データベース索引LSH検索+線形LSH検索との5種類が用意されている。

【0108】

厳密検索は、元のクエリベクトルをそのまま用いて、クエリベクトルに類似する事例ベクトルを線形検索する検索方式である。線形LSH検索は、LSH即値索引情報40Aを用いて、クエリハッシュに完全一致する事例ハッシュに対応する事例ベクトルを線形検索する検索方式である。データベース索引LSH検索は、LSH近傍展開索引情報40Bを用いて、クエリハッシュに完全一致する周辺ハッシュに対応する事例ベクトルを検索する検索方式である。データベース索引LSH検索+厳密検索は、データベース索引LSH検索により件数を絞り込んだ後に厳密検索を行う2段階絞り込み方式の検索方式である。データベース索引LSH検索+線形LSH検索は、データベース索引LSH検索により件数を絞

10

20

30

40

50

り込んだ後に線形 L S H 検索を行う 2 段階絞込方式の検索方式である。これら 5 種類の方式による検索は、それぞれ、図 1 3 に示した厳密検索器 3 3 1、線形 L S H 検索器 3 3 2、データベース索引 L S H 検索器 3 3 3、データベース索引 L S H 検索 + 厳密検索器 3 3 4 およびデータベース索引 L S H 検索 + 線形 L S H 検索器 3 3 5 により行われる。

【 0 1 0 9 】

また、本実施形態では、ベクトル類似性判定処理部 3 3 0 が、事前に与えられた検索方式設定情報 9 0 ( 図 1 3 参照 ) に基づいて、上記の 5 種類の検索方式のいずれかを検索対象のレコード数に応じて選択するものとする。検索方式設定情報 9 0 は、以下のような複数行からなる文法となっている。件数条件が上から照合され、マッチした検索方式が採用される。

件数条件 : 検索方式

件数条件 : 検索方式

...

otherwise : 検索方式

【 0 1 1 0 】

件数条件には、“ <=NNNN ” ( ただし NNNN は任意の正数 ) または “ otherwise ” を記述することができる。“ <=NNNN ” は検索対象となる事例ベクトルの数が NNNN 件以下の場合にマッチする。“ otherwise ” はあらゆる件数にマッチする。“ otherwise ” は検索方式設定情報 9 0 内に必ず 1 行含まれている必要がある。

【 0 1 1 1 】

検索方式には、“ Strict ” ( 厳密検索 )、 “ LinearLSH ” ( 線形 L S H 検索 )、 “ DBIndexedLSH ” ( データベース索引 L S H 検索 )、 “ DBIndexedLSH:NNNN/Strict ” ( データベース索引 L S H 検索 + 厳密検索 ) および “ DBIndexedLSH:NNNN/LinearLSH ” ( データベース索引 L S H 検索 + 線形 L S H 検索 ) の 5 種類を記述することができる。NNNN には任意の正数を記述することができる。

【 0 1 1 2 】

検索方式設定情報 9 0 の記述例を以下に示す。以下に例示する検索方式設定情報 9 0 は、検索対象となる事例ベクトルの数が 1 0 0 0 0 件以下なら厳密検索方式が採用され、 1 0 0 0 0 0 件以下なら線形 L S H 検索が採用され、それより多い件数ならデータベース索引 L S H 検索 + 線形 L S H 検索方式が採用されることを示している。

<=10000 : Strict

<=100000 : LinearLSH

Otherwise : DBIndexedLSH:100000/LinearLSH

【 0 1 1 3 】

図 1 5 は、ベクトル類似性判定部 3 3 0 による処理手順の一例を示すフローチャートである。ベクトル類似性判定部 3 3 0 は、検索対象データセット 7 0 を入力として受け取ると、例えば以下のステップ S 5 0 1 ~ ステップ S 5 0 7 の処理を実行して、検索結果 8 0 を出力する。

【 0 1 1 4 】

ステップ S 5 0 1 : ベクトル類似性判定部 3 3 0 は、入力された検索対象データセット 7 0 のレコード件数 ( 検索対象となる事例ベクトルの数 ) をカウントする。

【 0 1 1 5 】

ステップ S 5 0 2 : ベクトル類似性判定部 3 3 0 は、検索方式設定情報 9 0 を参照し、件数条件がステップ S 5 0 1 でカウントしたレコード件数とマッチする検索方式を取得する。取得した検索方式が Strict であればステップ S 5 0 3 に進み、 LinearLSH であればステップ S 5 0 4 に進み、 DBIndexedLSH であればステップ S 5 0 5 に進み、 DBIndexedLSH:NNN/Strict ( ただし NNNN は任意の正数 ) であればステップ S 5 0 6 に進み、 DBIndexedLSH:NNNN/LinearLSH ( ただし NNNN は任意の正数 ) であればステップ S 5 0 7 に進む。

【 0 1 1 6 】

ステップ S 5 0 3 : ベクトル類似性判定部 3 3 0 は、厳密検索器 3 3 1 を用いて検索結

10

20

30

40

50

果 8 0 を取得し、取得した検索結果 8 0 を出力する。

【 0 1 1 7 】

ステップ S 5 0 4 : ベクトル類似性判定部 3 3 0 は、線形 L S H 検索器 3 3 2 を用いて検索結果 8 0 を取得し、取得した検索結果 8 0 を出力する。

【 0 1 1 8 】

ステップ S 5 0 5 : ベクトル類似性判定部 3 3 0 は、データベース索引 L S H 検索器 3 3 3 を用いて検索結果 8 0 を取得し、取得した検索結果 8 0 を出力する。

【 0 1 1 9 】

ステップ S 5 0 6 : ベクトル類似性判定部 3 3 0 は、データベース索引 L S H 検索 + 厳密検索器 3 3 4 を用いて検索結果 8 0 を取得し、取得した検索結果 8 0 を出力する。

10

【 0 1 2 0 】

ステップ S 5 0 7 : ベクトル類似性判定部 3 3 0 は、データベース索引 L S H 検索 + 線形 L S H 検索器 3 3 5 を用いて検索結果 8 0 を取得し、取得した検索結果 8 0 を出力する。

【 0 1 2 1 】

なお、以上の例では検索対象となる事例ベクトルの数に応じて検索方式を切り替えるようにしているが、切り替えるべき検索方式の指定として、ビット長など様々なパラメタ設定も含ませるといった拡張を行うようにしてもよい。例えば 1 6 ビットなど短いビット長 (すなわち 1 件当たりの処理時間は小さいが低精度) の線形 L S H 検索で事例ベクトルの数を 1 0 0 0 0 0 件に絞り込んだ後、4 0 9 6 ビットなど長いビット長 (すなわち 1 件当

20

たりの処理時間は大きいが高精度) の線形 L S H 検索で順位付けするなどのプランを作ることができる。その場合の検索方式設定情報 9 0 の記述例を以下に示す (@以降にパラメタ設定が列挙される)。

<=100000 : LinearLSH@4096

otherwise : LinearLSH@16:100000/LinearLSH@4096

【 0 1 2 2 】

以下では、それぞれの検索方式による検索処理の具体例を説明する。

【 0 1 2 3 】

図 1 6 は、厳密検索器 3 3 1 による処理手順の一例を示すフローチャートである。厳密検索器 3 3 1 は、例えば以下のステップ S 6 0 1 ~ ステップ S 6 0 8 の処理を実行して、検索結果 8 0 を出力する。

30

【 0 1 2 4 】

ステップ S 6 0 1 : 厳密検索器 3 3 1 は、空の検索結果リストを生成する。

【 0 1 2 5 】

ステップ S 6 0 2 : 厳密検索器 3 3 1 は、vnn関数の第 2 引数として指定されたクエリベクトルを得る。

【 0 1 2 6 】

ステップ S 6 0 3 : 厳密検索器 3 3 1 は、検索対象データセット 7 0 からレコードを順に取り出す。

【 0 1 2 7 】

40

ステップ S 6 0 4 : 厳密検索器 3 3 1 は、ステップ S 6 0 3 で取り出したレコードにおいてvnn関数の第 1 引数として指定された事例ベクトルを格納する列に格納されている事例ベクトル群とクエリベクトルとのユークリッド距離を算出する。

【 0 1 2 8 】

ステップ S 6 0 5 : 厳密検索器 3 3 1 は、ステップ S 6 0 3 で取り出したレコードの行 ID とステップ S 6 0 4 で算出したユークリッド距離との組を検索結果リストに追加する。ただし、検索結果リストの要素がユークリッド距離昇順で並ぶような位置に挿入する。

【 0 1 2 9 】

ステップ S 6 0 6 : 厳密検索器 3 3 1 は、検索結果リストに含まれる要素数がvnn関数の第 3 引数として指定された出力上位件数よりも多い場合、出力上位件数以降の要素群を

50

破棄する。

【0130】

ステップS607：厳密検索器331は、検索対象データセット70からレコードをすべて取り出したか否かを判定する。そして、判定の結果がNoであればステップS603に戻って以降の処理を繰り返し、YesであればステップS608に進む。

【0131】

ステップS608：厳密検索器331は、検索結果リストに含まれる行IDの集合を検索結果80として出力し、処理を終了する。

【0132】

図17は、線形LSH検索器332による処理手順の一例を示すフローチャートである。線形LSH検索器332は、例えば以下のステップS701～ステップS709の処理を実行して、検索結果80を出力する。

【0133】

ステップS701：線形LSH検索器332は、空の検索結果リストを生成する。

【0134】

ステップS702：線形LSH検索器332は、vnn関数の引数として指定されたクエリベクトル、LSHアルゴリズムのアルゴリズム名および出力ビット長を取得して、クエリベクトルのハッシュ（縮約クエリベクトル）を生成する。以下、これを「クエリハッシュ」と呼ぶ。

【0135】

ステップS703：線形LSH検索器332は、検索対象データセット70のテーブル名、vnn関数の引数として指定された事例ベクトルを格納する列、LSHアルゴリズムのアルゴリズム名および出力ビット長を取得して、これらを元に決定されるテーブル名を持つLSH即値テーブルを取得する。

【0136】

ステップS704：線形LSH検索器332は、ステップS703で取得したLSH即値テーブルを用いて、検索対象データセット70の各レコードに対応する事例ハッシュを順に取り出す。

【0137】

ステップS705：線形LSH検索器332は、ステップS704で取り出した事例ハッシュとステップS702で生成したクエリハッシュとのハミング距離を算出する。

【0138】

ステップS706：線形LSH検索器332は、ステップS704で取り出した事例ハッシュに対応する検索対象データセット70のレコードの行IDとステップS705で算出したハミング距離との組を検索結果リストに追加する。ただし、検索結果リストの要素がハミング距離昇順で並ぶような位置に挿入する。

【0139】

ステップS707：線形LSH検索器332は、検索結果リストに含まれる要素数がvnn関数の第3引数として指定された出力上位件数よりも多い場合、出力上位件数以降の要素群を破棄する。

【0140】

ステップS708：線形LSH検索器332は、検索対象データセット70の各レコードに対応する事例ハッシュをすべて取り出したか否かを判定する。そして、判定の結果がNoであればステップS704に戻って以降の処理を繰り返し、YesであればステップS709に進む。

【0141】

ステップS709：線形LSH検索器332は、検索結果リストに含まれる行IDの集合を検索結果80として出力し、処理を終了する。

【0142】

なお、上記のステップS702において生成するクエリハッシュは本質的には二値ベク

10

20

30

40

50

トルであるが、LSH即値テーブルでの保存形式に合わせて、二進数と見立て整数として扱う。これにより、記憶領域を節約できるほか、整数同士の比較となるため照合処理が高速になる。

【0143】

図18は、データベース索引LSH検索器333による処理手順の一例を示すフローチャートである。データベース索引LSH検索器333は、例えば以下のステップS801～ステップS810の処理を実行して、検索結果80を出力する。

【0144】

ステップS801：データベース索引LSH検索器333は、空の検索結果セットを生成する。

10

【0145】

ステップS802：データベース索引LSH検索器333は、vnn関数の引数として指定されたクエリベクトル、LSHアルゴリズムのアルゴリズム名および出力ビット長を取得して、クエリハッシュを生成する。

【0146】

ステップS803：データベース索引LSH検索器333は、検索対象データセット70のテーブル名、vnn関数の引数として指定された事例ベクトルを格納する列、LSHアルゴリズムのアルゴリズム名および出力ビット長を取得して、これらを元に決定されるテーブル名を持つLSH近傍展開テーブルを取得する。

【0147】

ステップS804：データベース索引LSH検索器333は、ハミング距離条件変数に0を割り当てる。

20

【0148】

ステップS805：データベース索引LSH検索器333は、ステップS803で取得したLSH近傍展開テーブルに対して、HammingDistance列の値がハミング距離条件変数の現在値と一致し、かつ、HashValue列の値がステップS802で生成したクエリハッシュと一致するレコード群を検索し、得られた行IDのセットを検索結果セットに追加する。

【0149】

ステップS806：データベース索引LSH検索器333は、検索結果セットから重複要素を除去する。

30

【0150】

ステップS807：データベース索引LSH検索器333は、検索結果セットの要素数がvnn関数の第3引数として指定された出力上位件数以上となっているか否かを判定する。そして、判定の結果がNoであればステップS808に進み、YesであればステップS810に進む。

【0151】

ステップS808：データベース索引LSH検索器333は、ハミング距離条件変数の値を1追加する。

【0152】

ステップS809：データベース索引LSH検索器333は、ハミング距離条件変数の値がLSH近傍展開テーブルの距離上限以下か否かを判定する。そして、判定の結果がYesであればステップS805に戻って以降の処理を繰り返し、NoであればステップS810に進む。

40

【0153】

ステップS810：データベース索引LSH検索器333は、検索結果セットを検索結果80として出力し、処理を終了する。

【0154】

なお、上記のステップS802において生成するクエリハッシュは本質的には二値ベクトルであるが、LSH近傍展開テーブルでの保存形式に合わせて、二進数と見立て整数と

50

して扱う。これにより、記憶領域を節約できるほか、整数同士の比較となるため照合処理が高速になる。また、LSH近傍展開テーブルの距離上限が100などの大きな値の場合は、上記のステップS808での増分値を10など大きめの値とし、ステップS805でのハミング距離条件に範囲を設けるようにしてもよい。こうすれば反復回数が減るため、検索処理全体を高速化できる。

**【0155】**

図19は、データベース索引LSH検索+厳密検索器334による処理手順の一例を示すフローチャートである。データベース索引LSH検索+厳密検索は、データベース索引LSH検索により検索を行った後、その結果について厳密検索による検索を行う2段階絞込方式である。データベース索引LSH検索+厳密検索器334は、例えば以下のステップS901～ステップS904の処理を実行して、検索結果80を出力する。

10

**【0156】**

ステップS901：データベース索引LSH検索+厳密検索器334は、検索方式設定情報90に記載された検索方式名「DBIndexedLSH:NNNN/Strict」のNNNN部を読み取る。これが、データベース索引LSH検索における出力上位件数目標である。

**【0157】**

ステップS902：データベース索引LSH検索+厳密検索器334は、出力上位件数目標の設定でデータベース索引LSH検索器333を呼び出し、検索結果セットを得る。

**【0158】**

ステップS903：データベース索引LSH検索+厳密検索器334は、データベース索引LSH検索器333から取得した検索結果セットを検索対象データセット70として厳密検索器331を呼び出し、検索結果リストを得る。

20

**【0159】**

ステップS904：データベース索引LSH検索+厳密検索器334は、厳密検索器331から取得した検索結果リストに含まれる行IDの集合を検索結果80として出力し、処理を終了する。

**【0160】**

図20は、データベース索引LSH検索+線形LSH検索器335による処理手順の一例を示すフローチャートである。データベース索引LSH検索+線形LSH検索は、データベース索引LSH検索により検索を行った後、その結果について線形LSH検索による検索を行う2段階絞込方式である。データベース索引LSH検索+線形LSH検索器335は、例えば以下のステップS1001～ステップS1004の処理を実行して、検索結果80を出力する。

30

**【0161】**

ステップS1001：データベース索引LSH検索+線形LSH検索器335は、検索方式設定情報90に記載された検索方式名「DBIndexedLSH:NNNN/Strict」のNNNN部を読み取る。これが、データベース索引LSH検索における出力上位件数目標である。

**【0162】**

ステップS1002：データベース索引LSH検索+線形LSH検索器335は、出力上位件数目標の設定でデータベース索引LSH検索器333を呼び出し、検索結果セットを得る。

40

**【0163】**

ステップS1003：データベース索引LSH検索+線形LSH検索器335は、データベース索引LSH検索器333から取得した検索結果セットを検索対象データセット70として線形LSH検索器332を呼び出し、検索結果リストを得る。

**【0164】**

ステップS1004：データベース索引LSH検索+線形LSH検索器335は、線形LSH検索器332から取得した検索結果リストに含まれる行IDの集合を検索結果80として出力し、処理を終了する。

**【0165】**

50

ここで、上述した5種類の検索方式の選択基準について説明する。データ件数をN、特徴ベクトルの次元数をDとすると、厳密検索、線形LSH検索、データベース索引LSH検索の3つ検索方式の計算コスト(積算相当回数)は以下となる。

厳密検索： $N \times D$

線形LSH検索： $N$  (LSHアルゴリズムの出力ビット長を64以下とした場合)

データベース索引LSH検索：(データベース索引コスト)

ここで、 $N$ よりはるかに小さいため、上記3つの検索方式の中ではデータベース索引LSH検索が最も低コストである。

【0166】

精度については、厳密検索が最も高く、次に線形LSH検索が高く、データベース索引LSH検索が最も低い。近似的手法であるLSHを用いない厳密検索が最も高精度なのは当然だが、仮に同じLSHアルゴリズムを用いていても、線形LSH検索よりもデータベース索引LSH検索のほうが精度が低くなる理由は、データベース索引LSH検索では大きなビット長、ハミング距離を扱えないためである。データベース索引LSH検索は、ハミング距離以下のすべてのハッシュをレコードとして展開するLSH近傍展開テーブルを内部参照する。ハミング距離をH、ビット長をLとすると、レコード数Sは下記式(1)で表される。

【数1】

$$S = N \sum_h^H L C_h \quad \dots(1)$$

【0167】

この式(1)から分かるように、大きなハミング距離、ビット長のLSHを扱うと、レコード数が爆発するという問題がある。このため、小さなハミング距離、ビット長のLSHを扱うことになるが、その場合、距離分解能が低下し、実際に出力する件数を要求された出力件数に正確に合わせることが困難となる。この問題は、線形LSH検索、データベース索引LSH検索など、LSHを用いるいずれの検索手法でも起こりえるが、特に小さなハミング距離、ビット長のLSHしか扱えないデータベース索引LSH検索において顕著である。この欠点を解決するには、データベース索引LSH検索+厳密検索や、データベース索引LSH検索+線形LSH検索のように、出力件数を制御しやすい検索方式を2段階目に持つ2段階絞込方式を用いるようにし、1段階目のデータベース索引LSH検索では多めの出力件数で出力し、2段階目の検索で出力件数を要求に正確に合わせるようにすることが有効である。

【0168】

以上のことから、検索速度を許容範囲内に収めつつなるべく高い精度を得るためには、検索対象データセット70の件数が小規模の場合は厳密検索、検索対象データセット70の件数が中規模の場合は線形LSH検索、検索対象データセット70の件数が大規模の場合は、データベース索引LSH検索+厳密検索、あるいはデータベース索引LSH検索+線形LSH検索を用いることが有効である。

【0169】

以上、具体的な例を挙げながら詳細に説明したように、本実施形態のデータ管理システムでは、事前処理によってLSH近傍展開索引情報40Bを含む索引情報40を生成し、検索時にはこの索引情報40を用いて、クエリベクトルに完全一致する周辺ベクトルに対応する事例ベクトルを特定するといった計算量の少ない方法でデータテーブル20に対するベクトル近傍検索を行うようにしている。したがって、従来の類似検索において計算量の多くを占めていたベクトル近傍検索の計算量を削減して、データ検索の実行時間を短縮することができる。

【0170】

<第2実施形態>

次に、第2実施形態のデータ管理システムについて説明する。本実施形態は、検索時に

10

20

30

40

50

クエリベクトルに対する距離が所定値以下の周辺ベクトルを生成し、周辺ベクトルと完全一致する事例ベクトルを見つける構成とすることで、メモリ使用量の削減およびそれに伴うLSHの精度向上を実現したシステムである。

【0171】

上述した第1実施形態では、データベース索引LSH検索を実現するために、内部でLSH近傍展開テーブルやデータベース複合索引を用いている。これらの構造データは、上述したように、ハッシュのビット長およびハミング距離上限が大きいとサイズが肥大化する。しかし、データベース400の記憶領域として用いるディスクなどの記憶容量には限界があるため、大きなビット長およびハミング距離上限は扱えない場合がある。そして、十分に大きなビット長やハミング距離が扱えない場合、検索精度の低下を起す可能性がある。

10

【0172】

本実施形態では、第1実施形態のように検索対象となる各事例ベクトルのハッシュ（二値ベクトル）である事例ハッシュの近傍にあるハッシュ（二値ベクトル）である周辺ハッシュ群を列挙する代わりに、クエリベクトルのハッシュ（二値ベクトル）の近傍にあるハッシュ（二値ベクトル）群を周辺ハッシュ群として列挙する方式を取る。この方式の場合、1つのベクトルに対して列挙したハッシュ群を記憶するだけでよいため、記憶容量の限界という問題は解消され、検索精度の低下を起すリスクが回避される。ただし、第1実施形態と比較してクエリが複雑となるため、第1実施形態よりも検索速度は低下する。

【0173】

20

本実施形態のデータ管理システムにおける基本的な枠組みは第1実施形態と同様である。ただし、本実施形態では、検索時にLSH近傍展開索引情報40Bを用いないため、索引構築器200（図5参照）は、LSH即値索引情報40Aを生成するLSH即値索引情報生成器220のみを備えた構成となる。また、検索時にクエリハッシュを起点とした周辺ハッシュを生成するため、検索器300のベクトル類似性判定処理部330における検索方式が第1実施形態とは異なる。以下では、第1実施形態からの主な変更点となる本実施形態の検索方式について説明する。

【0174】

図21は、本実施形態の検索器300Aの構成例を示すブロック図である。第1実施形態との違いは、ベクトル類似性判定処理部330Aが、図13に示したデータベース索引LSH検索器333、データベース索引LSH検索+厳密検索器334およびデータベース索引LSH検索+線形LSH検索器335に代えて、クエリ摂動型LSH検索器336、クエリ摂動型LSH検索+厳密検索器337およびクエリ摂動型LSH検索+線形LSH検索器338を備える点である。

30

【0175】

クエリ摂動型LSH検索+厳密検索器337およびクエリ摂動型LSH検索+線形LSH検索器338については、2段階絞込方式の前段の検索方式がデータベース索引LSH検索からクエリ摂動型LSH検索に代わる他は第1実施形態のデータベース索引LSH検索+厳密検索器334およびデータベース索引LSH検索+線形LSH検索器335と同様である。このため、以下では、クエリ摂動型LSH検索+厳密検索器337およびクエリ摂動型LSH検索+線形LSH検索器338の説明は省略し、クエリ摂動型検索を行うクエリ摂動型LSH検索器336についてのみ説明する。

40

【0176】

図22は、クエリ摂動型LSH検索器336による処理手順の一例を示すフローチャートである。クエリ摂動型LSH検索器336は、例えば以下のステップS1101～ステップS1112の処理を実行して、検索結果80を出力する。

【0177】

ステップS1101：クエリ摂動型LSH検索器336は、空の検索結果セットを生成する。

【0178】

50

ステップ S 1 1 0 2 : クエリ摂動型 L S H 検索器 3 3 6 は、vnn関数の引数として指定されたクエリベクトル、L S Hアルゴリズムのアルゴリズム名および出力ビット長を取得して、クエリハッシュを生成する。

【 0 1 7 9 】

ステップ S 1 1 0 3 : クエリ摂動型 L S H 検索器 3 3 6 は、検索対象データセット 7 0 のテーブル名、vnn関数の引数として指定された事例ベクトルを格納する列、L S Hアルゴリズムのアルゴリズム名および出力ビット長を取得して、これらを元に決定されるテーブル名を持つ L S H 即値テーブルを取得する。

【 0 1 8 0 】

ステップ S 1 1 0 4 : クエリ摂動型 L S H 検索器 3 3 6 は、ハミング距離条件変数に 0 を割り当てる。

【 0 1 8 1 】

ステップ S 1 1 0 5 : クエリ摂動型 L S H 検索器 3 3 6 は、ステップ S 1 1 0 2 で生成したクエリハッシュからの距離がハミング距離条件変数の現在地と一致する周辺ハッシュをすべて生成する。

【 0 1 8 2 】

ステップ S 1 1 0 6 : クエリ摂動型 L S H 検索器 3 3 6 は、ステップ S 1 1 0 5 で生成した周辺ハッシュを順に取り出す。

【 0 1 8 3 】

ステップ S 1 1 0 7 : クエリ摂動型 L S H 検索器 3 3 6 は、ステップ S 1 1 0 3 で取得した L S H 即値テーブルから、ステップ S 1 1 0 6 で取り出した周辺ハッシュと同じ値を持つレコード群をすべて取り出し、各レコードの行 I D を検索結果セットに追加する。

【 0 1 8 4 】

ステップ S 1 1 0 8 : クエリ摂動型 L S H 検索器 3 3 6 は、ステップ S 1 1 0 5 で生成した周辺ハッシュをすべて取り出したか否かを判定する。そして、判定の結果が N o であればステップ S 1 1 0 6 に戻って以降の処理を繰り返し、判定の結果が Y e s であればステップ S 1 1 0 9 に進む。

【 0 1 8 5 】

ステップ S 1 1 0 9 : クエリ摂動型 L S H 検索器 3 3 6 は、検索結果セットから重複要素を除去する。

【 0 1 8 6 】

ステップ S 1 1 1 0 : クエリ摂動型 L S H 検索器 3 3 6 は、検索結果セットの要素数が vnn関数の第 3 引数として指定された出力上位件数以上となっているか否かを判定する。そして、判定の結果が N o であればステップ S 1 1 1 1 に進み、Y e s であればステップ S 1 1 1 2 に進む。

【 0 1 8 7 】

ステップ S 1 1 1 1 : クエリ摂動型 L S H 検索器 3 3 6 は、ハミング距離条件変数の値を 1 追加する。

【 0 1 8 8 】

ステップ S 1 1 1 2 : クエリ摂動型 L S H 検索器 3 3 6 は、検索結果セットを検索結果 8 0 として出力し、処理を終了する。

【 0 1 8 9 】

以上説明したように、本実施形態のデータ管理システムでは、クエリベクトルに類似する周辺ベクトルを生成し、この周辺ベクトルに完全一致する事例ベクトルを特定するといった計算量の少ない方法でデータテーブル 2 0 に対するベクトル近傍検索を行うようにしている。したがって、第 1 実施形態と同様に、従来の類似検索において計算量の多くを占めていたベクトル近傍検索の計算量を削減して、データ検索の実行時間を短縮することができる。

【 0 1 9 0 】

また、本実施形態では、第 1 実施形態と比較して記憶容量に対する制約が緩和されるこ

10

20

30

40

50

とで、比較的大きなビット長およびハミング距離上限を扱うことが可能になるため、ビット長およびハミング距離上限の制限に伴って検索精度の低下を起こすといったリスクが回避される。

【 0 1 9 1 】

< 第 3 実施形態 >

次に、第 3 実施形態のデータ管理システムについて説明する。本実施形態は、ベクトルの写像に用いる L S H の方式が第 1 実施形態とは異なる。すなわち、第 1 実施形態では写像処理にビットワイズ L S H を用いていたが、本実施形態では、写像処理に直積量子化 L S H を用いる。

【 0 1 9 2 】

直積量子化 L S H は、与えられたベクトルを整数ベクトル（整数の要素からのみ構成される多値ベクトル）に変換する。この整数ベクトルに対して第 1 実施形態で示したような索引手法を導入することで、検索時におけるベクトル距離計算を排除することができる。なお、第 1 実施形態と同様に、取り扱う事例ベクトルおよびクエリベクトルが元々低次元の整数ベクトルであるならば、直積量子化 L S H による写像は不要であり、後述の事例ハッシュとして事例ベクトルをそのまま用い、クエリハッシュとしてクエリベクトルをそのまま用いればよい。また、この場合、後述の周辺ハッシュとして、事例ベクトルを元に生成される周辺ベクトルを用いればよい。

【 0 1 9 3 】

直積量子化 L S H の代表的なアルゴリズムは、下記の参考文献 4 に記載されている。直積量子化 L S H としては、この参考文献 4 に記載のアルゴリズムを含め、様々なアルゴリズムが提案されている。

参考文献 4 : Herve Jegou , Matthijs Douze and Cordelia Schmid , “ Product quantization for nearest neighbor search ” , Pattern Analysis and Machine Intelligence , IEEE Transactions on 33.1 ( 2011 ) : 117-128 .

【 0 1 9 4 】

直積量子化 L S H によるハッシュ生成の代表的な手順は以下の通りである。

- ( 1 ) まず、与えられた事例ベクトル群が属する空間から、（好ましくは互いに排他的な）部分空間を複数取得する。
- ( 2 ) 各部分空間において、事例ベクトル群を予め定めた数のクラスタに分類する。
- ( 3 ) 各事例ベクトルについて、各部分空間のクラスタ番号を並べた整数ベクトルを生成する。これが写像後のベクトルである。
- ( 4 ) 各クラスタの分布などを用いて、写像後のベクトル間の差異から、写像前のベクトルの距離を近似的に算出するための距離モデルを生成する。

【 0 1 9 5 】

上記の手順（ 4 ）で生成される距離モデルはアルゴリズムによって異なるが、いずれの距離モデルであっても、ハッシュである 2 つの整数ベクトルが与えられたときに、元のベクトル空間での距離を概算できる性質を持つ。また多くのアルゴリズムでは、クラスタ番号が隣接する場合、元のベクトル空間におけるクラスタの分布領域も近接している。次元数は変えず、元のベクトル空間において格子状に並んだ最近傍の点で代表させる格子ベクトル量子化は、最も単純な直積量子化 L S H アルゴリズムである。

【 0 1 9 6 】

本実施形態のデータ管理システムにおける基本的な枠組みは第 1 実施形態と同様である。本実施形態の第 1 実施形態からの変更点は、整数ベクトルをハッシュとして扱うこと、および、ベクトル間距離としてハミング距離を用いる代わりに直積量子化 L S H が提供する距離モデルによって概算できる距離を用いるようにすることである。以下では、これらの変更点について説明する。

【 0 1 9 7 】

図 2 3 は、本実施形態の索引構築器 2 0 0 B の構成例を示すブロック図である。本実施形態の索引構築器 2 0 0 B は、図 2 3 に示すように、第 1 実施形態の索引構築器 2 0 0 に

10

20

30

40

50

おける L S H 即値索引情報生成器 2 2 0 ( 図 5 参照 ) に代えて P Q L S H 即値索引情報生成器 2 4 0 を備えるとともに、第 1 実施形態の索引構築器 2 0 0 における L S H 近傍展開索引情報生成器 2 3 0 ( 図 5 参照 ) に代えて P Q L S H 近傍展開索引情報生成器 2 5 0 を備える。

#### 【 0 1 9 8 】

P Q L S H 即値索引情報生成器 2 4 0 は、索引構築命令タイプ分類器 2 1 0 から索引構築命令文 3 0 を受け取った場合に、この索引構築命令文 3 0 に従って、P Q L S H 即値テーブルと、この P Q L S H 即値テーブルの HashValue 列に対するデータベース索引とを含む P Q L S H 即値索引情報 4 0 C を生成する。P Q L S H 即値索引情報生成器 2 4 0 が生成する P Q L S H 即値索引情報 4 0 C は、P Q L S H 即値テーブルの HashValue 列に格納される値が二値ベクトルではなく直積量子化 L S H アルゴリズムを用いて生成された整数ベクトルとなる点を除き、第 1 実施形態で説明した L S H 即値索引情報 4 0 A と同様である。このため、P Q L S H 即値索引情報生成器 2 4 0 については詳細な説明を省略する。

10

#### 【 0 1 9 9 】

P Q L S H 近傍展開索引情報生成器 2 5 0 は、索引構築命令タイプ分類器 2 1 0 から索引構築命令文 3 0 を受け取った場合に、この索引構築命令文 3 0 に従って、P Q L S H 近傍展開テーブルと、この P Q L S H 近傍展開テーブルの HashValue 列および Distance 列に対するデータベース複合索引 ( コンポジットインデックス ) とを含む P Q L S H 近傍展開索引情報 4 0 D を生成する。

#### 【 0 2 0 0 】

20

図 2 4 は、P Q L S H 近傍展開テーブルの一例を示す図である。P Q L S H 近傍展開テーブルは、図 2 4 に示すように、HashValue 列 ( 型 : blob )、Distance 列 ( 型 : real ) および DataRowID 列 ( 型 : integer ) からなるテーブルデータである。この図 2 4 に示す例では、要素値 { 0, 1, 2 } を取る長さ 2 の整数ベクトルのハッシュを用い、距離上限を 1 . 0 とした。Distance 列の値は、直積量子化 L S H アルゴリズムの提供する距離モデルによって概算される距離である。本例では、ハッシュの第 1 要素値が 1 変化すると距離が 0 . 4、第 2 要素値が 1 変化すると距離が 1 . 0 変化するという距離モデルを用いるものとした。

#### 【 0 2 0 1 】

P Q L S H 近傍展開索引情報生成器 2 5 0 は、索引構築命令文 3 0 で指定されたデータテーブル 2 0 から特徴ベクトルを取り出し、各特徴ベクトル ( 事例ベクトル ) について、索引構築命令文 3 0 で指定された L S H アルゴリズムおよび距離モデルを用いて整数ベクトルである事例ハッシュを生成する。次に、P Q L S H 近傍展開索引情報生成器 2 5 0 は、生成した各事例ハッシュについて、索引構築命令文 3 0 で指定された距離以下の周辺ハッシュ群を生成する。そして、P Q L S H 近傍展開索引情報生成器 2 5 0 は、例えば、周辺ハッシュと、周辺ハッシュの事例ハッシュからの距離と、事例ハッシュの元になる事例ベクトルが属するデータテーブル 2 0 のレコードの行 ID との 3 つ組からなるレコードを、P Q L S H 近傍展開テーブルに登録する。

30

#### 【 0 2 0 2 】

なお、P Q L S H 近傍展開テーブルの HashValue 列および Distance 列に対するデータベース複合索引は、図 1 0 に示した第 1 実施形態のデータベース複合索引とほとんど同じであり、図 1 0 の HammingDistance の値に沿った分岐が、Distance の値に沿った分岐に変わるだけである。

40

#### 【 0 2 0 3 】

図 2 5 は、P Q L S H 近傍展開索引情報生成器 2 5 0 による処理手順の一例を示すフローチャートである。P Q L S H 近傍展開索引情報生成器 2 5 0 は、索引構築命令タイプ分類器 2 1 0 から索引構築命令文 3 0 を受け取ると、例えば以下のステップ S 1 2 0 1 ~ ステップ S 1 2 1 0 の処理を実行し、P Q L S H 近傍展開テーブルとデータベース複合索引とを含む P Q L S H 近傍展開索引情報 4 0 D をデータベース 4 0 0 に保存する。

#### 【 0 2 0 4 】

50

ステップ S 1 2 0 1 : P Q L S H 近傍展開索引情報生成器 2 5 0 は、データベース 4 0 0 内に、名前 : HashValue ・ 型 : blob の列と、名前 : Distance ・ 型 : real の列と、名前 : DataRowID ・ 型 : integer の列とを持つテーブル ( 空の P Q L S H 近傍展開テーブル ) を生成する。なお、テーブル名の形式は「AAAA\_BBBB\_CCCC\_DDDD」とする。ただし、AAAA は処理対象のデータテーブル 2 0 のテーブル名、BBBB は処理対象の列の列名、CCCC は直積量子化 L S H アルゴリズムのアルゴリズム名、DDDD はモデル ( アルゴリズムに関する設定や分析結果 ) 名とする。例えば、aTable\_feat\_PQ\_confA といったテーブル名が生成したテーブルに与えられる。

【 0 2 0 5 】

ステップ S 1 2 0 2 : P Q L S H 近傍展開索引情報生成器 2 5 0 は、索引構築命令文 3 0 で指定されたデータテーブル 2 0 から、索引構築命令文 3 0 で指定された列の特徴ベクトルを順に取り出す。

10

【 0 2 0 6 】

ステップ S 1 2 0 3 : P Q L S H 近傍展開索引情報生成器 2 5 0 は、ステップ S 1 2 0 2 で取り出した特徴ベクトル ( 事例ベクトル ) に対して、索引構築命令文 3 0 で指定された直積量子化 L S H アルゴリズムおよびモデルを用いて、事例ハッシュを生成する。

【 0 2 0 7 】

ステップ S 1 2 0 4 : P Q L S H 近傍展開索引情報生成器 2 5 0 は、ステップ S 1 2 0 3 で得られた事例ハッシュからの距離が指定距離以下の周辺ハッシュをすべて生成する。ただし、起点となる事例ハッシュと各周辺ハッシュとの距離は、直積量子化 L S H アルゴリズムのモデルを用いて概算する。

20

【 0 2 0 8 】

ステップ S 1 2 0 5 : P Q L S H 近傍展開索引情報生成器 2 5 0 は、ステップ S 1 2 0 4 で得られた周辺ハッシュを順に取り出す。

【 0 2 0 9 】

ステップ S 1 2 0 6 : P Q L S H 近傍展開索引情報生成器 2 5 0 は、ステップ S 1 2 0 5 で取り出した周辺ハッシュ、起点となった事例ハッシュからの距離、事例ハッシュの元になる事例ベクトルが属するデータテーブル 2 0 内のレコードの行 ID の 3 つ組からなるレコードを、P Q L S H 近傍展開テーブルに追加する。

【 0 2 1 0 】

30

ステップ S 1 2 0 7 : P Q L S H 近傍展開索引情報生成器 2 5 0 は、ステップ S 1 2 0 4 で得られた周辺ハッシュをすべて取り出したか否かを判定する。そして、判定の結果が N o であればステップ S 1 2 0 5 に戻って以降の処理を繰り返し、Y e s であればステップ S 1 2 0 8 に進む。

【 0 2 1 1 】

ステップ S 1 2 0 8 : P Q L S H 近傍展開索引情報生成器 2 5 0 は、処理対象のデータテーブル 2 0 から特徴ベクトルをすべて取り出したか否かを判定する。そして、判定の結果が N o であればステップ S 1 2 0 2 に戻って以降の処理を繰り返し、Y e s であればステップ S 1 2 0 9 に進む。

【 0 2 1 2 】

40

ステップ S 1 2 0 9 : P Q L S H 近傍展開索引情報生成器 2 5 0 は、以上の処理を経て生成された P Q L S H 近傍展開テーブルの HashValue 列および Distance 列に対するデータベース複合索引 ( コンポジットインデックス ) を構築し、データベース 4 0 0 に保存する。

【 0 2 1 3 】

ステップ S 1 2 1 0 : P Q L S H 近傍展開索引情報生成器 2 5 0 は、事例ハッシュの生成に用いた直積量子化 L S H アルゴリズムのモデルを、テーブル名と同じファイル名 ( 例えば aTable\_feat\_PQ\_confA ) でデータベース 4 0 0 に保存する。

【 0 2 1 4 】

なお、上記のステップ S 1 2 0 4 で周辺ハッシュを生成する方法は第 1 実施形態と似た

50

方法を用いるが、本実施形態では、ハッシュのどの要素を変更するかだけでなく、どれだけ値を変更するかも決定する必要がある。これには、例えばダイクストラ法などの最良優先探索などを用いればよい。

【0215】

図26は、本実施形態の検索器300Bの構成例を示すブロック図である。第1実施形態との違いは、ベクトル類似性判定処理部330Bが、図13に示した線形LSH検索器332、データベース索引LSH検索器333、データベース索引LSH検索+厳密検索器334およびデータベース索引LSH検索+線形LSH検索器335に代えて、線形PQLSH検索器351、データベース索引PQLSH検索器352、データベース索引PQLSH検索+厳密検索器353およびデータベース索引PQLSH検索+線形PQLSH検索器354を備える点である。

10

【0216】

線形PQLSH検索器351は、扱うベクトルが第1実施形態からの実装変更に伴って二値ベクトルから整数ベクトルに変わる他は、第1実施形態の線形LSH検索器332と同様である。また、データベース索引PQLSH検索+厳密検索器353およびデータベース索引PQLSH検索+線形PQLSH検索器354については、2段階絞込方式の前段の検索方式がデータベース索引LSH検索からデータベース索引PQLSH検索に代わる他は、第1実施形態のデータベース索引LSH検索+厳密検索器334およびデータベース索引LSH検索+線形LSH検索器335と同様である。このため、以下では、線形PQLSH検索器351、データベース索引PQLSH検索+厳密検索器353およびデータベース索引PQLSH検索+線形PQLSH検索器354の説明は省略し、データベース索引PQLSH検索を行うデータベース索引PQLSH検索器352についてのみ説明する。なお、線形LSH検索については、第1実施形態から変更せずにそのまま用いることも可能である。

20

【0217】

図27は、データベース索引PQLSH検索器352による処理手順の一例を示すフローチャートである。データベース索引PQLSH検索器352は、例えば以下のステップS1301～ステップS1310の処理を実行して、検索結果80を出力する。

【0218】

ステップS1301：データベース索引PQLSH検索器352は、空の検索結果セットを生成する。

30

【0219】

ステップS1302：データベース索引PQLSH検索器352は、vnn関数の引数として指定されたクエリベクトル、直積量子化LSHアルゴリズムのアルゴリズム名およびオプションパラメタを取得して、クエリハッシュを生成する。

【0220】

ステップS1303：データベース索引PQLSH検索器352は、検索対象データセット70のテーブル名、vnn関数の引数として指定された事例ベクトルを格納する列、直積量子化LSHアルゴリズムのアルゴリズム名およびオプションパラメタを取得して、これらを元に決定されるテーブル名を持つPQLSH近傍展開テーブルを取得する。

40

【0221】

ステップS1304：データベース索引PQLSH検索器352は、距離条件変数に0を割り当てる。

【0222】

ステップS1305：データベース索引PQLSH検索器352は、ステップS1303で取得したPQLSH近傍展開テーブルに対して、Distance列の値が距離条件変数の現在値以下であり、かつ、HashValue列の値がステップS1302で生成したクエリハッシュと一致するレコード群を検索し、得られた行IDのセットを検索結果セットに追加する。

【0223】

50

ステップ S 1 3 0 6 : データベース索引 P Q L S H 検索器 3 5 2 は、検索結果セットから重複要素を除去する。

【 0 2 2 4 】

ステップ S 1 3 0 7 : データベース索引 P Q L S H 検索器 3 5 2 は、検索結果セットの要素数が vnn 関数の第 3 引数として指定された出力上位件数以上となっているか否かを判定する。そして、判定の結果が N o であればステップ S 1 3 0 8 に進み、 Y e s であればステップ S 1 3 1 0 に進む。

【 0 2 2 5 】

ステップ S 1 3 0 8 : データベース索引 P Q L S H 検索器 3 5 2 は、距離条件変数を一定量追加する。

10

【 0 2 2 6 】

ステップ S 1 3 0 9 : データベース索引 P Q L S H 検索器 3 5 2 は、ハミング距離条件変数の値が P Q L S H 近傍展開テーブルの距離上限以下か否かを判定する。そして、判定の結果が Y e s であればステップ S 1 3 0 5 に戻って以降の処理を繰り返し、 N o であればステップ S 1 3 1 0 に進む。

【 0 2 2 7 】

ステップ S 1 3 1 0 : データベース索引 P Q L S H 検索器 3 5 2 は、検索結果セットを検索結果 8 0 として出力し、処理を終了する。

【 0 2 2 8 】

以上、具体的な例を挙げながら詳細に説明したように、本実施形態のデータ管理システムでは、事前処理によって P Q L S H 近傍展開索引情報 4 0 D を含む索引情報 4 0 を生成し、検索時にはこの索引情報 4 0 を用いて、クエリベクトルに完全一致する周辺ベクトルに対応する事例ベクトルを特定するといった計算量の少ない方法でデータテーブル 2 0 に対するベクトル近傍検索を行うようにしている。したがって、第 1 実施形態と同様に、従来の類似検索において計算量の多くを占めていたベクトル近傍検索の計算量を削減して、データ検索の実行時間を短縮することができる。

20

【 0 2 2 9 】

< 補足説明 >

上述した実施形態のデータ管理システムは、一例として、一般的なコンピュータとしてのハードウェアを用いた実行環境で動作するプログラムによる実装が可能である。この場合、本実施形態のデータ管理システムにおける上述の各機能的な構成要素（データ登録器 1 0 0、索引構築器 2 0 0（索引構築器 2 0 0 B）、検索器 3 0 0（検索器 3 0 0 A、検索器 3 0 0 B））は、ハードウェアとソフトウェア（プログラム）との協働により実現される。また、データベース 4 0 0 は、プログラムによってアクセス可能な任意のメモリ資源によって実現される。

30

【 0 2 3 0 】

図 2 8 は、データ管理システムのハードウェア構成例を示すブロック図である。データ管理システムは、例えば図 2 8 に示すように、CPU（Central Processing Unit）1 0 0 1 などのプロセッサ回路、ROM（Read Only Memory）1 0 0 2 や RAM（Random Access Memory）1 0 0 3 などの記憶装置、表示パネルや各種操作デバイスが接続される入出力 I / F 1 0 0 4、ネットワークに接続して通信を行う通信 I / F 1 0 0 5、各部を接続するバス 1 0 0 6 などを備えた、通常のコンピュータを利用したハードウェア構成とすることができる。

40

【 0 2 3 1 】

また、上述した構成のハードウェア上で実行されるプログラムは、例えば、インストール可能な形式または実行可能な形式のファイルで CD - ROM（Compact Disk Read Only Memory）、フレキシブルディスク（FD）、CD - R（Compact Disk Recordable）、DVD（Digital Versatile Disc）などのコンピュータで読み取り可能な記録媒体に記録されてコンピュータプログラムプロダクトとして提供される。また、上述した構成のハードウェア上で実行されるプログラムを、インターネットなどのネットワークに接続

50

されたコンピュータ上に格納し、ネットワーク経由でダウンロードさせることにより提供するように構成してもよい。また、上述した構成のハードウェア上で実行されるプログラムをインターネットなどのネットワーク経由で提供または配布するように構成してもよい。また、上述した構成のハードウェア上で実行されるプログラムを、ROM 1002などに予め組み込んで提供するように構成してもよい。

【0232】

上述した構成のハードウェア上で実行されるプログラムは、実施形態のデータ管理システムの各機能的な構成要素を含むモジュール構成となっており、例えば、CPU 1001（プロセッサ回路）が上記記録媒体からプログラムを読み出して実行することにより、上述した各部がRAM 1003（主記憶）上にロードされ、RAM 1003（主記憶）上に生成されるようになっている。なお、実施形態のデータ管理システムの各機能的な構成要素やデータベース400は、複数のコンピュータに跨って実現される構成であってもよい。また、上述の機能的な構成要素の一部または全部を、ASIC（Application Specific Integrated Circuit）やFPGA（Field-Programmable Gate Array）などの専用のハードウェアを用いて実現することも可能である。

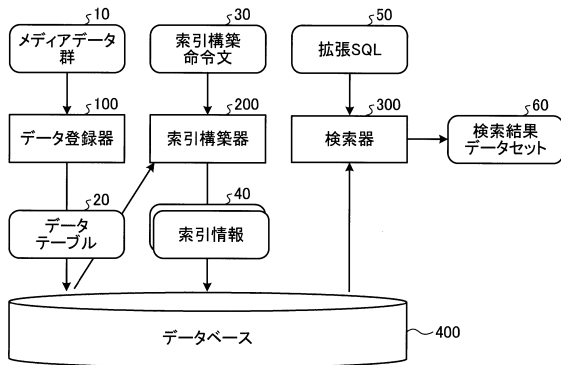
10

【0233】

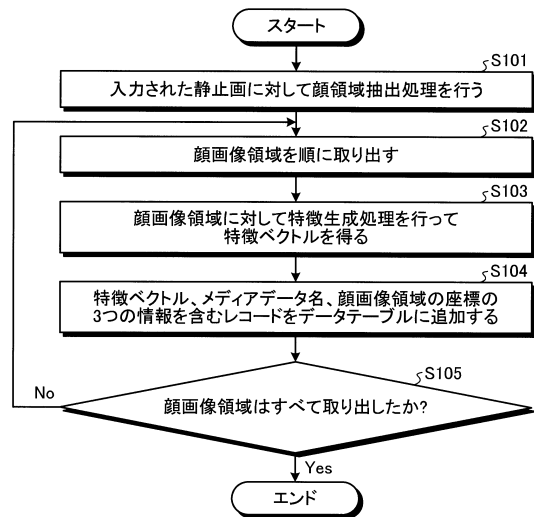
以上、本発明の実施形態を説明したが、ここで説明した実施形態は、例として提示したものであり、発明の範囲を限定することは意図していない。ここで説明した新規な実施形態は、その他の様々な形態で実施されることが可能であり、発明の要旨を逸脱しない範囲で、種々の省略、置き換え、変更を行うことができる。ここで説明した実施形態やその変形は、発明の範囲や要旨に含まれるとともに、請求の範囲に記載された発明とその均等の範囲に含まれる。

20

【図1】



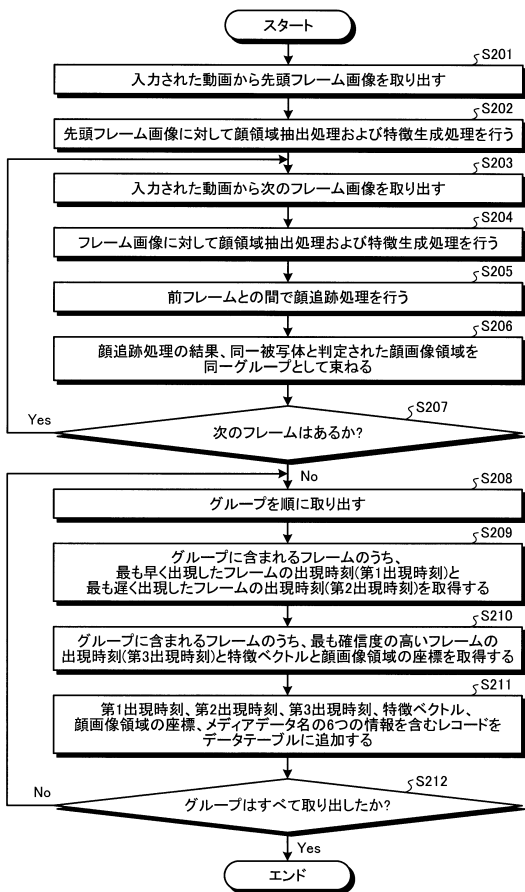
【図3】



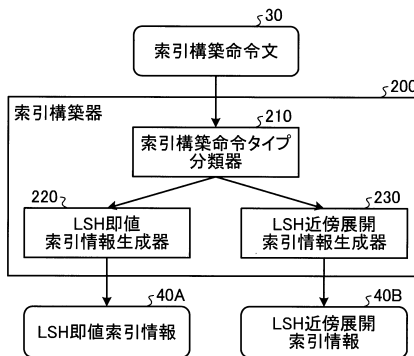
【図2】

file	feat	region	begin_time	end_time	best_time
dir/a.jpg	(32,44,33)	(10,10,300,300)			
dir/a.jpg	(76,11,74)	(300,300,600,600)			
dir/b.jpg	(76,12,64)	(400,400,800,800)			
dir/a.mpg	(10,33,42)	(250,280,450,650)	22:04:10:29	22:05:00:12	22:04:30:28
dir/a.mpg	(80,93,12)	(80,80,200,200)	32:24:21:12	33:04:40:31	32:28:22:29
⋮	⋮	⋮	⋮	⋮	⋮

【図4】



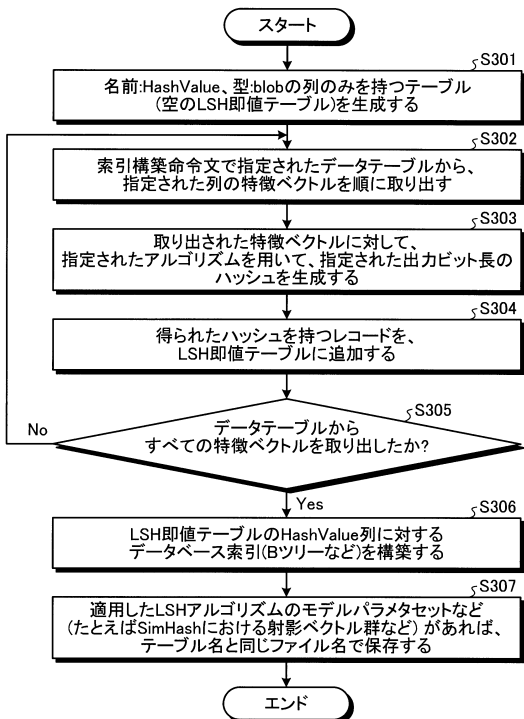
【図5】



【図6】

HashValue
100
010
000
001
011
⋮

【図7】



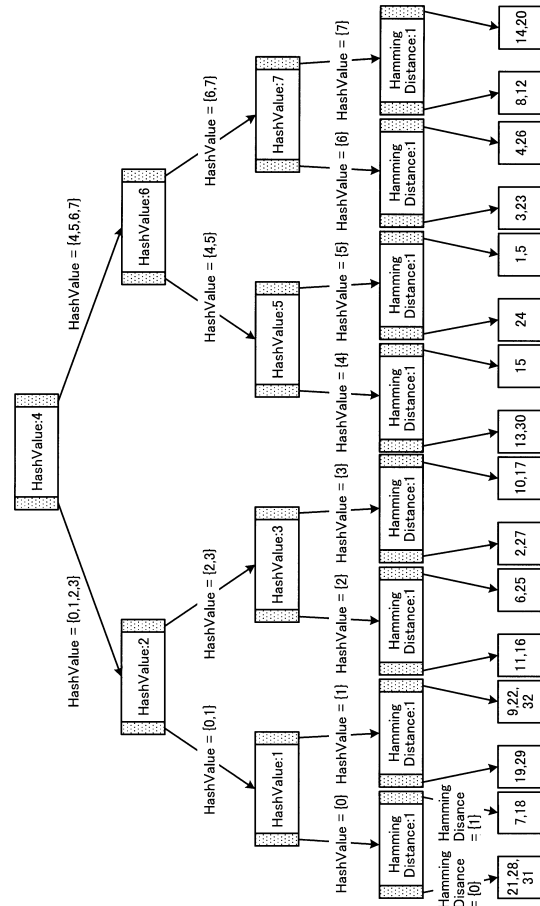
【図8】

HashValue	HammingDistance	DataRowID
100	0	1
000	1	1
110	1	1
101	1	1
010	0	2
110	1	2
000	1	2
011	1	2
000	0	3
100	1	3
010	1	3
001	1	3
001	0	4
101	1	4
011	1	4
000	1	4
011	0	5
111	1	5
001	1	5
010	1	5
⋮	⋮	⋮

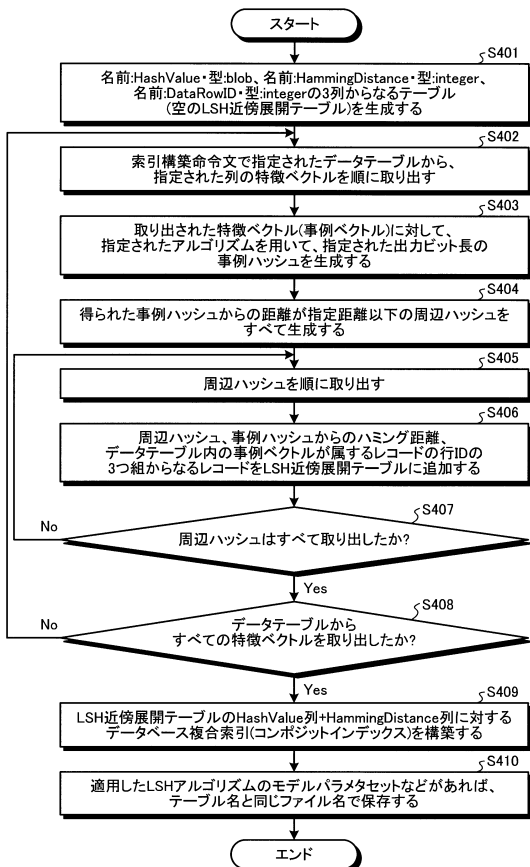
【図9】

(a)		(b)	
HashValue	HammingDistance	Table1_RowID	DataRowID
000	0	1	3
000	1	2	1
001	0	2	2
001	1	2	4
010	0	3	4
010	1	4	3
011	0	4	5
011	1	5	2
100	0	6	3
100	1	6	5
101	0	7	5
101	1	8	2
110	0	8	4
110	1	9	1
111	0	10	3
111	1	12	1
111	1	12	4
111	1	14	1
111	1	14	2
111	1	16	5
⋮	⋮	⋮	⋮

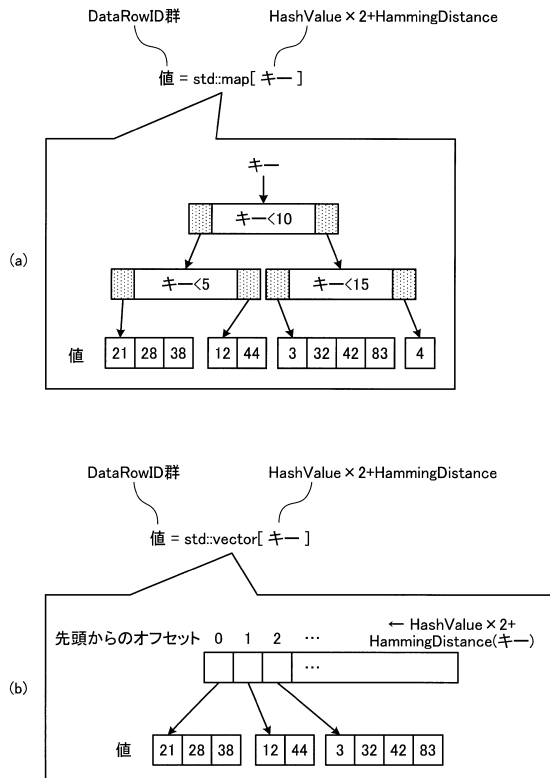
【図10】



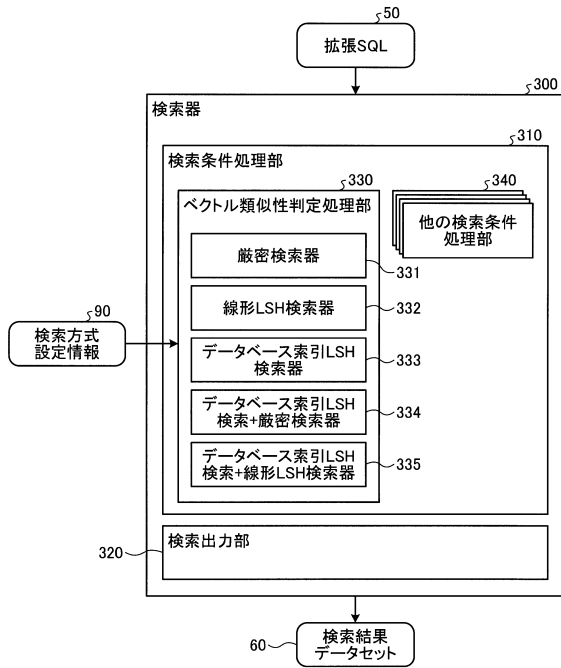
【図11】



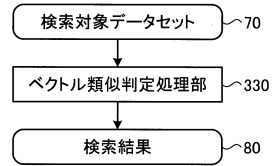
【図12】



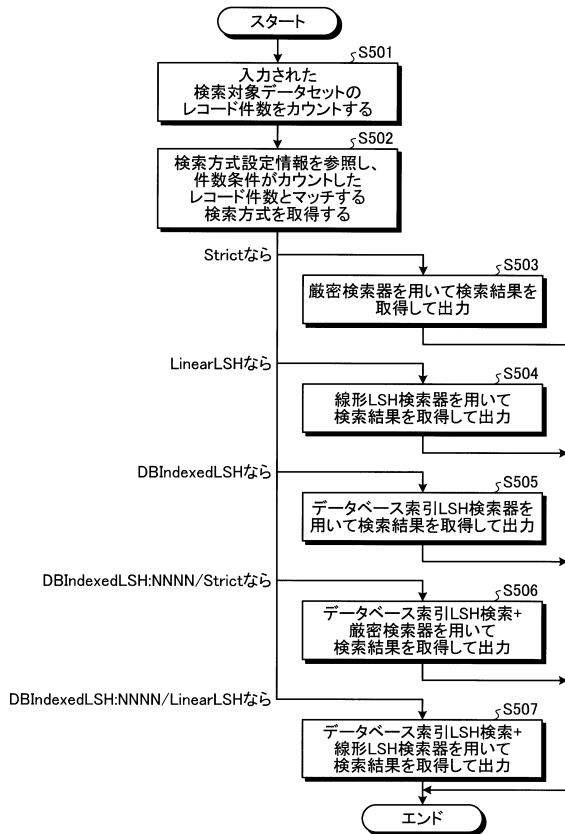
【図13】



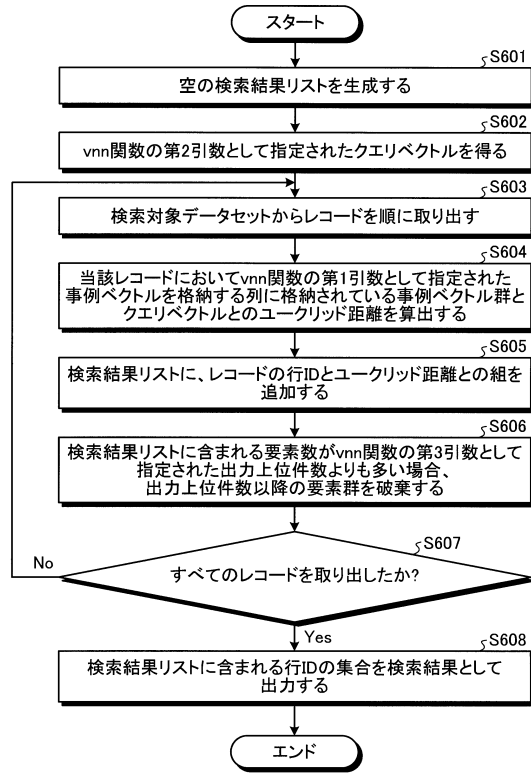
【図14】



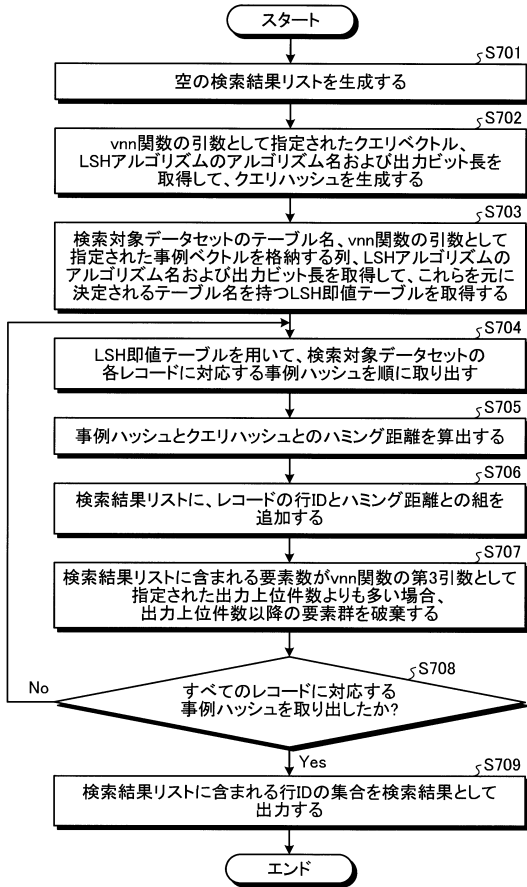
【図15】



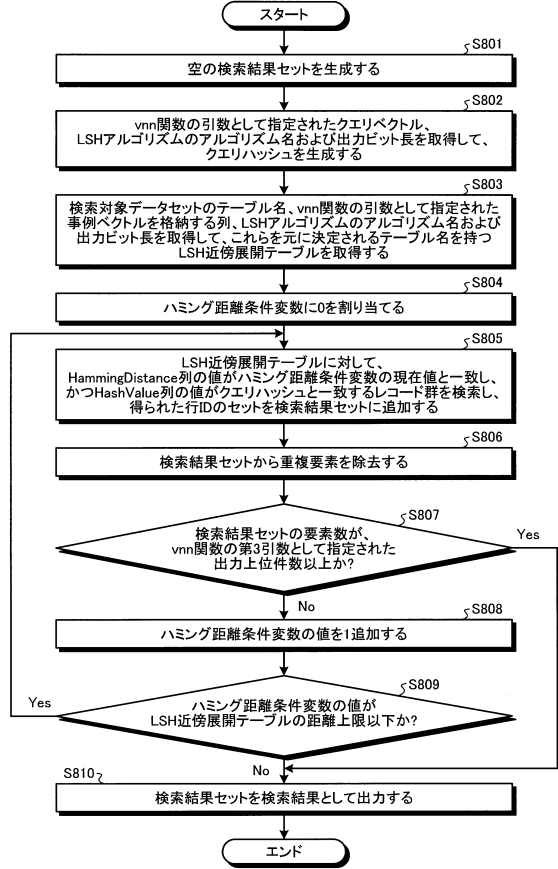
【図16】



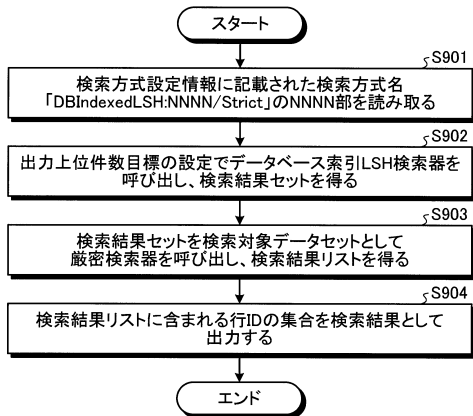
【図17】



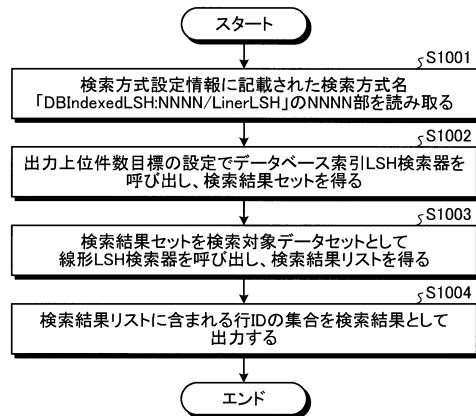
【図18】



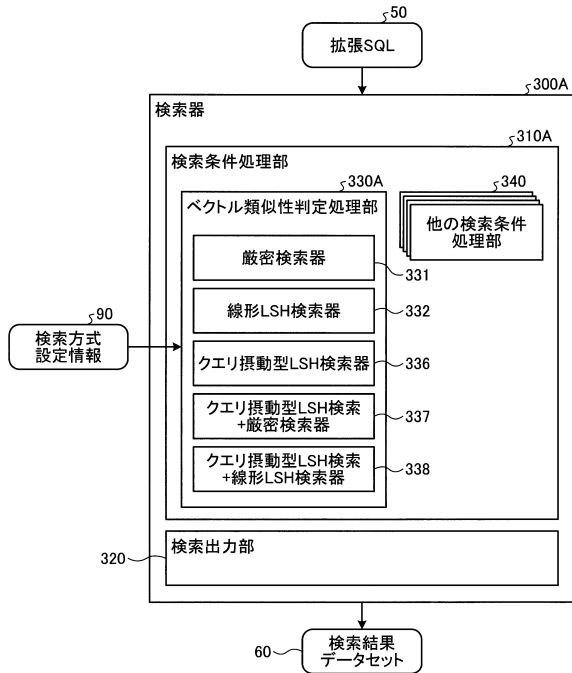
【図19】



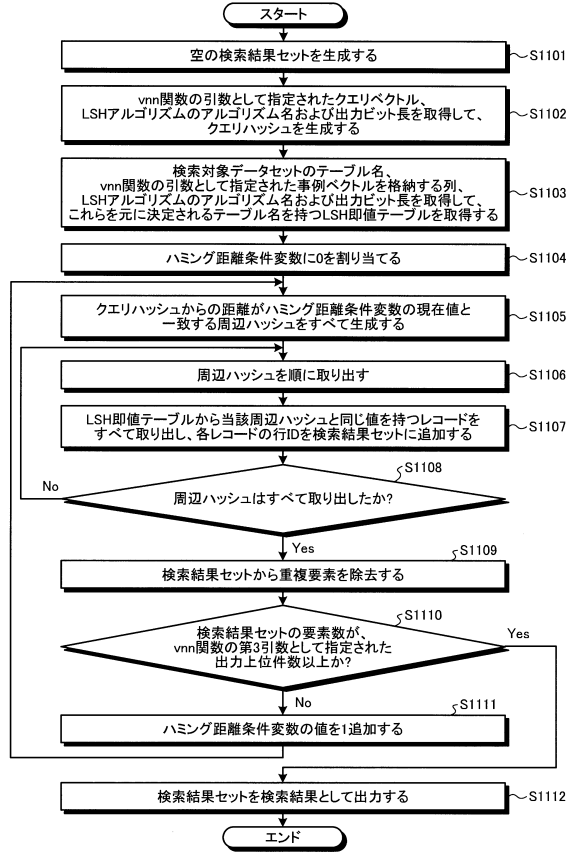
【図20】



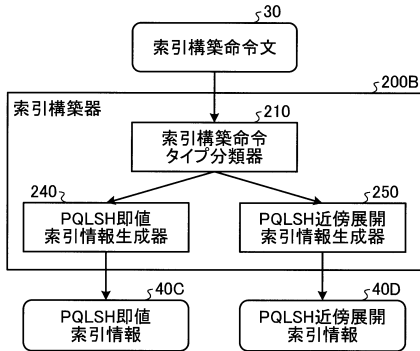
【図 2 1】



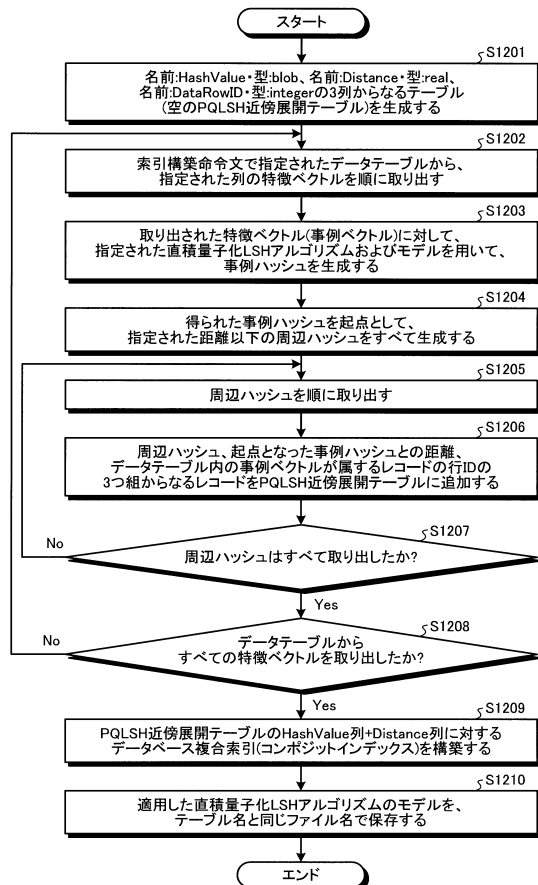
【図 2 2】



【図 2 3】



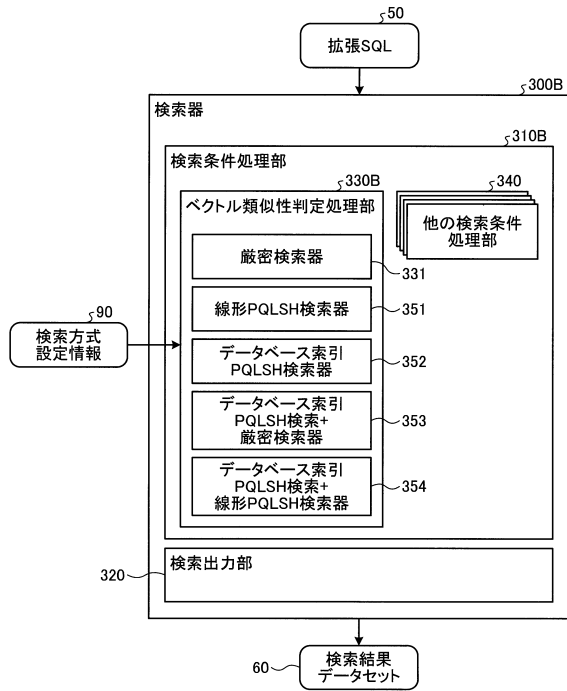
【図 2 5】



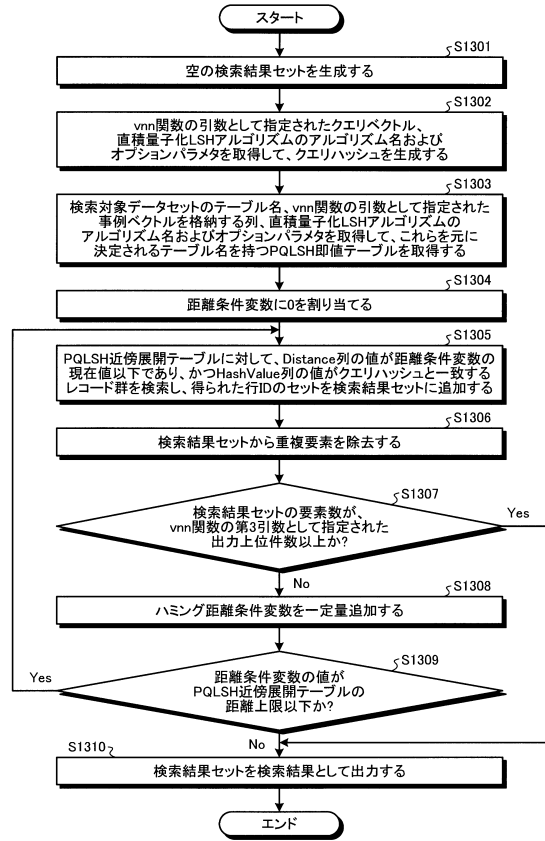
【図 2 4】

HashValue	Distance	DataRowID
10	0	1
00	0.4	1
20	0.4	1
11	1	1
20	0	2
10	0.4	2
00	0.8	2
21	1	2
⋮	⋮	⋮

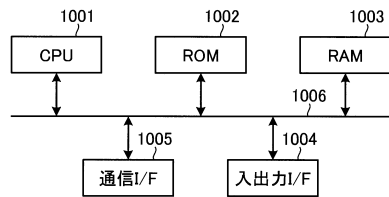
【図26】



【図27】



【図28】



---

フロントページの続き

- (72)発明者 湯浅 真由美  
東京都港区芝浦一丁目1番1号 株式会社東芝内
- (72)発明者 長田 邦男  
東京都港区芝浦一丁目1番1号 株式会社東芝内

審査官 早川 学

- (56)参考文献 国際公開第2013/129580(WO, A1)  
国際公開第2008/026414(WO, A1)  
BABENKO, Artem et al., The Inverted Multi-Index, 2012 IEEE Conference on Computer Vision and Pattern Recognition, IEEE, 2012年 6月21日, pp.3069-3076
- (58)調査した分野(Int.Cl., DB名)  
G06F 17/30