

(19) 日本国特許庁(JP)

(12) 公開特許公報(A)

(11) 特許出願公開番号

特開2005-149493

(P2005-149493A)

(43) 公開日 平成17年6月9日(2005.6.9)

(51) Int.Cl.⁷

G06F 17/30

F I

G06F 17/30

220B

テーマコード (参考)

5B075

G06F 17/30

170B

G06F 17/30

210C

審査請求 未請求 請求項の数 26 O L (全 39 頁)

(21) 出願番号 特願2004-314846 (P2004-314846)
 (22) 出願日 平成16年10月28日 (2004.10.28)
 (31) 優先権主張番号 60/515713
 (32) 優先日 平成15年10月31日 (2003.10.31)
 (33) 優先権主張国 米国 (US)
 (31) 優先権主張番号 10/729915
 (32) 優先日 平成15年12月9日 (2003.12.9)
 (33) 優先権主張国 米国 (US)

(71) 出願人 000005496
 富士ゼロックス株式会社
 東京都港区赤坂二丁目17番22号
 (74) 代理人 100079049
 弁理士 中島 淳
 (74) 代理人 100084995
 弁理士 加藤 和詳
 (72) 発明者 マシュー エル. クーパー
 アメリカ合衆国 94114 カリフォル
 ニア州 サンフランシスコ トウェンティ
 ーサード ストリート 3998
 (72) 発明者 ジョナサン ティー. フート
 アメリカ合衆国 94025 カリフォル
 ニア州 メンロ パーク ローレル スト
 リート 450

最終頁に続く

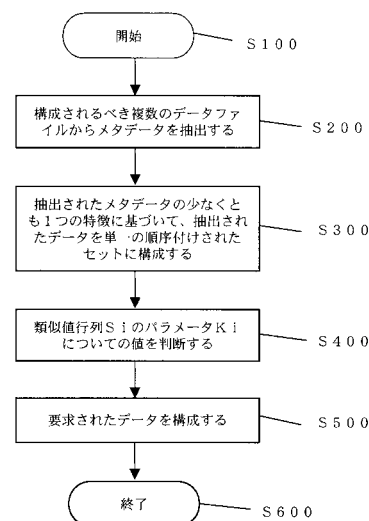
(54) 【発明の名称】 データファイルを構成するための方法、プログラム、及びシステム

(57) 【要約】

【課題】複数のデータファイルを構成するための方法、プログラム、及びシステムを提供する。

【解決手段】データ構成システム及び方法は、データファイルの少なくともいくつかについて、関連したデータを抽出することによって、複数のデータファイルに関連するメタデータ又は他のデータを用いて、複数のデータファイルを構成し、抽出された関連したデータ及び入力パラメータ値に基づいて、抽出された関連したデータに構成して、データファイルを少なくともいくつかの群に分割する。

【選択図】 図1



【特許請求の範囲】

【請求項 1】

少なくとも各データファイルと関連した少なくとも 1 つのメタデータ要素を有するメタデータを用いて、複数のデータファイルを構成するための方法であって、

前記データファイルの少なくともいくつかについて、前記データファイルと関連した少なくとも 1 つのメタデータ要素を抽出し、

前記抽出されたメタデータ要素についての値に基づいて、前記抽出されたメタデータ要素を要求された順序に構成し、

少なくとも 1 つのパラメータ値を入力し、

前記抽出されたメタデータ要素及び前記入力パラメータ値に基づいて、前記データファイルの少なくともいくつかを群に分割する、ことを含む方法。 10

【請求項 2】

前記少なくともいくつかのデータファイルを分割することが、前記少なくとも 1 つのパラメータ値の少なくとも 1 つの各々について、前記抽出されたメタデータ要素及びそのパラメータ値を用いて、前記複数のデータファイルの少なくとも 2 つについての類似性値を判断することを含む、請求項 1 に記載の方法。

【請求項 3】

前記少なくとも 1 つの類似性値を判断することが、次の数式で少なくとも 1 つの類似性値を判断することを含み、 20

$$S_K(i, j) = \exp\left(-\frac{|t_i - t_j|}{K}\right)$$

$S_K(i, j)$ が i^{th} データファイル及び j^{th} データファイルについての類似性値であり、

K がパラメータ値であり、

t_i 及び t_j が、前記 i^{th} 及び j^{th} データファイルについての前記少なくとも 1 つの抽出されたメタデータ要素の少なくとも一つのメタデータ要素の実際の値である、 30

請求項 2 に記載の方法。

【請求項 4】

少なくとも 1 つの類似性値を判断することが、次の数式で前記少なくとも 1 つの類似性値を判断することを含み、

$$S_K(i, j) = \exp\left(\frac{1}{K} \left(\frac{\langle v_i, v_j \rangle}{|v_i| |v_j|} - 1 \right) \right)$$

$S_K(i, j)$ が i^{th} データファイル及び j^{th} データファイルについての類似性値であり、

K がパラメータ値であり、

v_i 及び v_j が、前記 i^{th} 及び j^{th} データファイルから判断された実際のベクトル値である、

請求項 2 に記載の方法。

【請求項 5】

少なくともいくつかのデータファイルの各々について、そのデータファイルについての及び多数の近くのデータファイルについての少なくとも 1 つの類似性値に基づいて、その 50

データファイルについての少なくとも 1 つの新規性値を判断することをさらに含む、請求項 2 に記載の方法。

【請求項 6】

少なくとも 1 つの新規性値を判断することが、次の数式で少なくとも 1 つの新規性値を判断することを含み、

$$v_K(s) = \sum_{l,n=-5}^5 S_K(s+l, s+n)g(l,n)$$

10

$v_K(s)$ が新規性値であり、

g がガウシアンテーパの $1 \times 1 \times 1$ チェッカーボードカーネルである、

請求項 5 に記載の方法。

【請求項 7】

前記データファイルの少なくともいくつかについて判断された前記少なくとも 1 つの新規性値に基づいて、前記複数のデータファイルの境界位置間の少なくとも 1 つの境界位置を判断することをさらに含む、請求項 5 に記載の方法。

【請求項 8】

前記判断された境界位置の少なくともいくつかについて、その境界位置についての信頼値を判断することをさらに含む、請求項 7 に記載の方法。

20

【請求項 9】

境界位置についての信頼値を判断することが、次の数式で信頼値を判断することを含み、

$$C(B_K) = \sum_{l=1}^{|B_K|-1} \frac{1}{(b_{l+1} - b_l)^2} \sum_{i,j=b_l}^{b_{l+1}} S_K(i, j) - \sum_{l=1}^{|B_K|-2} \frac{1}{(b_{l+1} - b_l)(b_{l+2} - b_{l+1})} \sum_{i=b_l}^{b_{l+1}} \sum_{j=b_{l+1}}^{b_{l+2}} S_K(i, j)$$

$C(B_K)$ が B_K の境界についての信頼値であり、

$S_K(i, j)$ が i 番データファイル及び j 番データファイルについての類似性値であり、

30

b が前記入力パラメータ K レベルについての特定の値で検出された境界のインデックス値である、

請求項 8 に記載の方法。

【請求項 10】

前記判断された境界位置の少なくともいくつかについて、前記信頼値を最大にする少なくとも 1 つのパラメータ値の少なくとも 1 つを判断することをさらに含む、請求項 8 に記載の方法。

【請求項 11】

データファイルの対応するメタデータ要素と少なくとも関連した少なくとも 1 つのメタデータ要素を有するメタデータを用いて、複数のデータファイルを構成するための方法であって、

40

各メタデータがデータファイルに対応している、少なくとも 1 つのメタデータのセットを処理することと、

前記メタデータを分析するために要求されたパラメータ値を得ることと、

得られたパラメータ値を用いてメタデータ要素のセット内の構造を判断することであって、前記複数のデータファイルの少なくともサブセットについて、前記パラメータ値を互いに用いて前記メタデータの少なくともサブセットを比較することによって前記構造が判断される、構造を判断することと、を含む、

方法

50

【請求項 1 2】

前記メタデータ要素の前記判断された構造を用いて、前記データファイルを群にクラスタ化することをさらに含む、請求項 1 1 に記載の方法。

【請求項 1 3】

データファイルの前記判断されたクラスタから境界を判断することをさらに含み、前記境界はデータファイルの前記判断されたクラスタ間に位置している、請求項 1 2 に記載の方法

【請求項 1 4】

データファイルの 1 つのクラスタ内の前記メタデータ要素の少なくともいくつかを、データファイルの要素クラスタ内の前記メタデータ要素の少なくともいくつかと比較することによって、類似性値を判断することと、

データファイルの 1 つのクラスタ内の前記メタデータ要素の少なくともいくつかを、データファイルの他のクラスタ内の前記メタデータ要素の少なくともいくつかと比較することによって、非類似性値を判断することと、をさらに含む、

請求項 1 3 に記載の方法。

【請求項 1 5】

前記類似性値と前記非類似性値との相違に基づいて、データファイルのクラスタの要求された群に対応するパラメータ値を判断することをさらに含む、請求項 1 4 に記載の方法

【請求項 1 6】

データ処理装置上で実行可能であると共に、少なくとも各データファイルと関連した少なくとも 1 つのメタデータ要素を有するメタデータを用いることによって複数のデータファイルを構成するために用いられ得る、プログラムであって、前記プログラムが、

前記データファイルの少なくともいくつかについて、そのデータファイルと関連した少なくとも 1 つのメタデータ要素を抽出するための命令と、

前記抽出されたメタデータ要素についての値に基づいて、前記抽出されたメタデータ要素を要求された順序に構成するための命令と、

パラメータ値を入力するための命令と、

前記抽出されたメタデータ要素及び前記入力パラメータ値に基づいて、前記データファイルの少なくともいくつかを群に分割するための命令と、を含む、

プログラム。

【請求項 1 7】

前記データファイルの少なくともいくつかを群に分割するための命令が、前記少なくとも 1 つのパラメータ値の少なくとも 1 つの各々について、前記抽出されたメタデータ要素の少なくともいくつか及びそのパラメータ値を用いて、前記複数のデータファイルの少なくとも 2 つについての類似性値を判断するための命令をさらに含む、請求項 1 6 に記載のプログラム。

【請求項 1 8】

少なくともいくつかのデータファイルの各々について、そのデータファイルについての及び多数の近くのデータファイルについての少なくとも 1 つの類似性値に基づいて、そのデータファイルについての少なくとも 1 つの新規性値を判断するための命令をさらに含む、請求項 1 7 に記載のプログラム。

【請求項 1 9】

前記少なくとも 1 つの類似性値を判断するための命令が、次の数式で前記少なくとも 1 つの類似性値を判断するための命令を含み、

$$S_K(i, j) = \exp\left(-\frac{|t_i - t_j|}{K}\right)$$

10

20

30

40

50

$S_K(i, j)$ が i^{th} データファイル及び j^{th} データファイルについての類似性値であり、

K がパラメータ値であり、

t_i 及び t_j が、前記 i^{th} 及び j^{th} データファイルについての前記少なくとも 1 つの抽出されたメタデータ要素の少なくとも一つのメタデータ要素の実値である、

請求項 17 に記載のプログラム。

【請求項 20】

少なくとも 1 つの類似性値を判断するための命令が、次の数式で前記少なくとも 1 つの類似性値を判断するための命令を含み、

$$S_K(i, j) = \exp \left(\frac{1}{K} \left(\frac{\langle v_i, v_j \rangle}{\|v_i\| \|v_j\|} - 1 \right) \right)$$

10

$S_K(i, j)$ が i^{th} データファイル及び j^{th} データファイルについての類似性値であり、

K がパラメータ値であり、

v_i 及び v_j が、前記 i^{th} 及び j^{th} データファイルから判断された実際のベクトル値である、

20

請求項 17 に記載のプログラム。

【請求項 21】

前記データファイルの少なくともいくつかについて判断された前記少なくとも 1 つの新規性値に基づいて、前記複数のデータファイルの境界位置間の少なくとも 1 つの境界位置を判断するための命令をさらに含む、請求項 18 に記載のプログラム。

【請求項 22】

少なくとも 1 つの新規性値を判断するための命令が、次の数式で少なくとも 1 つの新規性値を判断する命令を含み、

$$v_K(s) = \sum_{l, n=-5}^5 S_K(s+l, s+n) g(l, n)$$

30

$v_K(s)$ が新規性値であり、

g がガウシアンテーパーの 11×11 チェッカーボードカーネルである、

請求項 18 に記載のプログラム。

【請求項 23】

前記判断された境界位置の少なくともいくつかについて、その境界位置についての信頼値を判断するための命令をさらに含む、請求項 21 に記載のプログラム。

【請求項 24】

前記少なくとも 1 つの信頼値を判断するための命令が、次の数式でこのような信頼値の各々を判断するための命令を含み、

40

$$C(B_K) = \sum_{l=1}^{|B_K|-1} \frac{1}{(b_{l+1} - b_l)^2} \sum_{i, j=b_l}^{b_{l+1}} S_K(i, j) - \sum_{l=1}^{|B_K|-2} \frac{1}{(b_{l+1} - b_l)(b_{l+2} - b_{l+1})} \sum_{i=b_l}^{b_{l+1}} \sum_{j=b_{l+1}}^{b_{l+2}} S_K(i, j)$$

$C(B_K)$ が B_K^{th} 境界についての信頼値であり、

$S_K(i, j)$ が前記 i^{th} データファイル及び前記 j^{th} データファイルについての類似性値であり、

50

bがあるレベルの検出された境界である、
請求項23に記載のプログラム。

【請求項25】

前記判断された境界位置の少なくともいくつかについて、前記信頼値を最大にする少なくとも1つのパラメータ値の少なくとも1つを判断するための命令をさらに含む、請求項23に記載のプログラム。

【請求項26】

データファイルの対応するメタデータ要素と少なくとも関連した、少なくとも1つのメタデータ要素を有するメタデータを用いて、複数のデータファイルを構成するために用いられ得る、データファイル構成システムであって、

前記データファイルの少なくともいくつかについて、前記データファイルと関連した少なくとも1つのメタデータ要素を抽出する、メタデータ抽出回路、ルーティン、又はアプリケーションと、

前記抽出されたメタデータについての値に基づいて、前記抽出されたメタデータ要素を要求された順序に構成するための、メタデータ構成回路、ルーティン、又はアプリケーションと、

前記少なくとも1つのパラメータ値の少なくとも1つについて、前記抽出されたメタデータ要素の少なくともいくつか及びそのパラメータ値を用いて、前記複数のデータファイルの少なくとも2つについての類似性値を判断する、類似性値判断回路、ルーティン、又はアプリケーションと、

前記データファイルについての及び多数の近くのデータファイルについての少なくとも1つの類似性値に基づいて、そのデータファイルについての少なくとも1つの新規性値を判断する、新規性値判断回路、ルーティン、又はアプリケーションと、

前記データファイルの少なくともいくつかについて判断された前記少なくとも1つの新規性値に基づいて前記複数のデータファイルの境界位置間の少なくとも1つの境界位置を判断することによって、前記抽出されたメタデータ要素及び前記入力パラメータ値に基づいて、前記データファイルの少なくともいくつかを群に分割する、データ分割回路、ルーティン、又はアプリケーションと、

前記判断された境界位置の少なくともいくつかについて、その境界位置についての信頼値を判断し、前記判断された境界位置の少なくともいくつかについて、前記データ分割回路、ルーティン、又はアプリケーションが、前記信頼値を最大にする前記少なくとも1つのパラメータ値をさらに判断する、信頼値判断回路、ルーティン、又はアプリケーションと、を備える、

データファイル構成システム。

【発明の詳細な説明】

【技術分野】

【0001】

本発明は、データの階層的クラスタリングによってデータを構成するためのシステム及び方法に関する。すなわち、複数のデータファイルを構成するための方法、プログラム、及びシステムに関する。

【背景技術】

【0002】

データは、例えば、メディアデータとして、メディアファイル内等に、種々の方法によって保存されている。メディアデータは、オーディオ、ビデオ、グラフィック、及び/又はテキストストリーム又はファイル等の、ストリーム又はファイルでもよい。メディアデータの1つの例示的な形式は、デジタル写真である。高品質のデジタルカメラの値ごろ感、デジタル写真を増加させ、多くの人々がデジタル写真を容易に撮影して保存できるようにした。これらのデジタル写真は、しばしば、デジタル写真データファイルとして保存される。

【0003】

メディアデータファイルは、通常、いくつかの異なる部分を含む。たとえば、1つのデジタル写真データファイルは、たとえばJ P E Gフォーマット等の特定のファイルフォーマットで記録された画像データを含んでいてもよい。画像データに加えて、画像データについてのある情報は、画像データと関連した得られるデジタル写真データファイル内に、メタデータとして典型的に保存されてもよい。関連したメタデータは、基礎となる画像データとは分離した、異なるデータである。1つの例示的なフォーマットは、Exif (エグジフ、Exchangeable Image File Format) であり、それは、J P E G画像データファイルの一部として保存されたヘッダ情報についてのフォーマットとして、しばしば用いられる。Exifフォーマットで保存されたメタデータの例は、ファイル名、データが作成された時間、画像ファイルに最後の変更が行なわれた時間等の1つ以上のタイムスタンプ、画像データの短い説明、又は画像データが得られた場所についてのG P S位置を含む。

10

【0004】

デジタル写真データファイル及び他のこのように急速に蓄積するデータファイルを処理するために多くの技術が作成されてきた。単純なデータファイルについては、このような技術の1つは、このようなデータファイルの各々が関連した項目に応じて、このようなデータファイルを特定のフォルダ内に配置することを含む。他の技術は、ある人の連絡先の情報を、パーソナルコンピュータデータベース内の所与のファイルディレクトリ内に手動で構成することを含む。ユーザは、内容を検討し、ファイルディレクト内への特定の連絡先の情報の配置、及び友人、仕事関係 (business contact)、学校関係 (school contact) 等の任意のサブカテゴリを判断する。

20

【0005】

Microsoft Word (登録商標) で用いられるフォーマット等特定のフォーマットで書かれた連絡先の情報のような、単純なデータでさえ、二つの特徴を含む。データを識別するデータレコードの名前は、レコード内に含まれた情報を圧縮するスカラー特徴と呼び得る。連絡先の名前、連絡先 (contact) の住所、又はその特定の連絡先に関する他のデータ等の記録の実際の内容は、より詳細であり、ベクトル特徴と呼び得る。

【0006】

【特許文献1】米国特許第6,542,869 B1号

【非特許文献1】ヘッカーマン (Heckerman), 「ベイズのネットワークを学ぶに当たっての指導書」 ("A Tutorial on Learning With Bayesian Networks"), マイクロソフトリサーチ (Microsoft Research), 1995年3月, P. 1 ~ 57

30

【非特許文献2】フット (Foote), 「オーディオの新規性の測定を用いる自動的オーディオ分割」 ("Automatic Audio Segmentation Using a Measure of Audio Novelty"), F X パロ・アルト研究所 (FX Palo Alto Laboratory, Inc)

【非特許文献3】プラットら (Platt et al.), 「写真T O C : 個人的な写真を閲覧するための自動的クラスタリング」 ("Photo TOC: Automatic Clustering for browsing Personal Photographs"), マイクロソフト・リサーチ (Microsoft Research), 2002年2月, P. 1 ~ 19

【非特許文献4】スラニーら (Slaney et al.), 「マルチメディア・エッジ: すべての次元における階層を見つける」 ("Multimedia Edges: Finding Hierarchy in all Dimensions"), マルチメディアについての第9回A C M国際会議の予稿集 (Proceedings of the 9th ACM International Conference on Multimedia), P. 1 ~ 12

40

【非特許文献5】ルイら (Loui et al.), 「アルバムに適用するための自動的な画像イベント分割及び品質選別」 ("Automatic Image Event Segmentation and Quality Screening for Albuming Applications"), マルチメディアについてのI E E E国際会議及び博覧会 (IEEE International Conference on Multimedia and Expo), 2000年7月, ニューヨーク

【非特許文献6】チェンら (Chen et al.), 「ベイズの情報基準による、話者、環境、及びチャンネル変化の検出及びクラスタリング」 ("Speaker, Environment and Channel Change Detection and Clustering Via the Bayesian Information Criterion"), I

50

B M T . J . ワトソン研究センター (IBM T.J.Watson Research Center)

【非特許文献 7】レナルら (Renals et al.) , 「会議室からのオーディオ情報アクセス」 (“ Audio Information Access From Meeting Rooms ”) , IEEE ICASSAP-2003 , 香港 , P . 1 ~ 4

【非特許文献 8】グラハムら (Graham et al.) , 「個人的なデジタルライブラリによる写真閲覧用の要素としての時間」 (“ Time as Essence for Photo Browsing Through Personal Digital Libraries ”) , スタンフォード大学 (Stanford University)

【非特許文献 9】「デジタルスチルカメラ画像ファイルフォーマット基準」 (“ Digital Still Camera Image File Format Standard ”) , Version2.1 , 1998 年 6 月 12 日、日本電子工業振興協会 (Japan Electronic Industry Development Association (JEIDA)) , P . 1 ~ 166

10

【発明の開示】

【発明が解決しようとする課題】

【0007】

データファイルを構成するための 1 つの方法は、ユーザが各データファイルの内容及び / 又はそのデータファイルの名前を実際に調査し、次に、適切な項目記述子が付けられたフォルダ等の、特定のファイルディレクトリ内の、データファイルの適切な位置を手動で判断することである。データファイルを特定の位置に配置して集めることは、データファイルを特定の関係に構成することである。しかしながら、たとえば、膨大な枚数の写真が編成されなければならない時には、写真の各データファイルを手動で編成・構成することは、ほとんど不可能になる。この困難さは、各データファイルの内容が複雑である場合、たとえば内容が画像データである場合に大きくなる。

20

【課題を解決するための手段】

【0008】

(本発明の目的)

本発明は、データファイル内のメタデータ又は他の順序付けられた情報に基づいてデータを効率良く構成するためのシステム及び方法を提供する。

【0009】

本発明は、データファイルを構成するメタデータに基づいて、関連したデータファイルをクラスタリングすることによって、データファイルを構成するための、システム及び方法を、別々に提供する。

30

【0010】

本発明は、データファイルのメタデータを抽出するためのシステム及び方法を、別々に提供する。

【0011】

本発明は、データファイルのメタデータに基づいてデータファイルを構成するためのシステム及び方法を、別々に提供する。

【0012】

本発明は、閲覧及び / 又は検索するために要求されたデータファイルを構成するためのシステム及び方法を、別々に提供する。

40

【0013】

本発明によるシステム及び方法の種々の例示的な実施の形態において、データファイルの要求されたセットは、メタデータのセットを調査することによって構成され、ここでは、メタデータの各メタデータ要素は、特定のデータファイルから抽出されるか、又は少なくとも関連して来た。種々の実施の形態において、メタデータのセット内の構造は、メタデータ要素を分析するためのメタデータの要素の値の要求された範囲を得て、次に、データファイルのすべて又はサブセットについてメタデータの要素についての値を比較することによって、評価される。

【0014】

種々の例示的な実施の形態において、メタデータのセットのメタデータ要素は、メタデ

50

ータのセットの評価された構造を用いて、クラスタ化される。メタデータのセットの構造は、メタデータ要素値の各クラスタを他のクラスタから線引きする境界を含む。種々の例示的な実施の形態において、比較されたデータファイル間の類似性又は非類似性を判断するために、あるデータファイルのあるメタデータ要素の値は、範囲の値に基づいて、クラスタ内の他のデータファイルのメタデータ要素の値と比較される。

【0015】

種々の例示的な実施の形態において、データは、データのすべての可能な対間の比較又はデータのすべての可能な対のサブセットを用いて、構成される。種々の例示的な実施の形態において、比較された類似性又は非類似性には、メタデータ要素のクラスタ及びそれらの対応するデータファイルの配置に対応する数値が与えられる。種々の例示的な実施の形態において、より正確にするために、クラスタの配置が調べられる。種々の例示的な実施の形態において、コンテンツベースの類似性測定を作り出すことによってデータファイルは低レベルの特徴を発生する時よりも、より効率良くかつコンピュータによってより安く構成される。

10

【0016】

本発明の第1の態様は、少なくとも各データファイルと関連した少なくとも1つのメタデータ要素を有するメタデータを用いて、複数のデータファイルを構成するための方法であって、データファイルの少なくともいくつかについて、データファイルと関連した少なくとも1つのメタデータ要素を抽出し、抽出されたメタデータ要素についての値に基づいて、抽出されたメタデータ要素を要求された順序に構成し、少なくとも1つのパラメータ値を入力し、抽出されたメタデータ要素及び入力パラメータ値に基づいて、データファイルの少なくともいくつかを群に分割する、ことを含む方法である。

20

【0017】

本発明の第2の態様は、第1の態様において、少なくともいくつかのデータファイルを分割することが、少なくとも1つのパラメータ値の少なくとも1つの各々について、抽出されたメタデータ要素及びそのパラメータ値を用いて、複数のデータファイルの少なくとも2つについての類似性値を判断することを含む、方法である。

【0018】

本発明の第3の態様は、第2の態様において、少なくとも1つの類似性値を判断することが、次の数式で少なくとも1つの類似性値を判断することを含み、

30

$$S_K(i, j) = \exp\left(-\frac{|t_i - t_j|}{K}\right)$$

$S_K(i, j)$ が i^{th} データファイル及び j^{th} データファイルについての類似性値であり、 K がパラメータ値であり、 t_i 及び t_j が、 i^{th} 及び j^{th} データファイルについての少なくとも1つの抽出されたメタデータ要素の少なくとも一つのメタデータ要素の実値である、方法である。

40

【0019】

本発明の第4の態様は、第2の態様において、少なくとも1つの類似性値を判断することが、次の数式で少なくとも1つの類似性値を判断することを含み、

$$S_K(i, j) = \exp\left(\frac{1}{K} \left(\frac{\langle v_i, v_j \rangle}{|v_i| |v_j|} - 1 \right) \right)$$

$S_K(i, j)$ が i^{th} データファイル及び j^{th} データファイルについての類似性値であ

50

り、 K がパラメータ値であり、 v_i 及び v_j が、 i^{th} 及び j^{th} データファイルから判断された実際のベクトル値である、方法である。

【0020】

本発明の第5の態様は、第2の態様において、少なくともいくつかのデータファイルの各々について、そのデータファイルについての及び多数の近くのデータファイルについての少なくとも1つの類似性値に基づいて、そのデータファイルについての少なくとも1つの新規性値を判断することをさらに含む、方法である。

【0021】

本発明の第6の態様は、第5の態様において、少なくとも1つの新規性値を判断することが、次の数式で少なくとも1つの新規性値を判断することを含み、

10

$$v_K(s) = \sum_{l,n=-5}^5 S_K(s+l, s+n) g(l, n)$$

$v_K(s)$ が新規性値であり、 g がガウシアンテーパの 11×11 チェッカーボードカーネルである、方法である。

【0022】

本発明の第7の態様は、第5の態様において、データファイルの少なくともいくつかについて判断された少なくとも1つの新規性値に基づいて、複数のデータファイルの境界位置間の少なくとも1つの境界位置を判断することをさらに含む、方法である。

20

【0023】

本発明の第8の態様は、第7の態様において、判断された境界位置の少なくともいくつかについて、その境界位置についての信頼値を判断することをさらに含む、方法である。

【0024】

本発明の第9の態様は、第8の態様において、境界位置についての信頼値を判断することが、次の数式で信頼値を判断することを含み、

$$C(B_K) = \sum_{l=1}^{|B_K|-1} \frac{1}{(b_{l+1} - b_l)^2} \sum_{i,j=b_l}^{b_{l+1}} S_K(i, j) - \sum_{l=1}^{|B_K|-2} \frac{1}{(b_{l+1} - b_l)(b_{l+2} - b_{l+1})} \sum_{i=b_l}^{b_{l+1}} \sum_{j=b_{l+1}}^{b_{l+2}} S_K(i, j)$$

30

$C(B_K)$ が B_K^{th} 境界についての信頼値であり、 $S_K(i, j)$ が i^{th} データファイル及び j^{th} データファイルについての類似性値であり、 b が入力パラメータ K レベルについての特定の値で検出された境界のインデックス値である、方法である。

【0025】

本発明の第10の態様は、第8の態様において、判断された境界位置の少なくともいくつかについて、信頼値を最大にする少なくとも1つのパラメータ値の少なくとも1つを判断することをさらに含む、方法である。

【0026】

本発明の第11の態様は、データファイルの対応するメタデータ要素と少なくとも関連した少なくとも1つのメタデータ要素を有するメタデータを用いて、複数のデータファイルを構成するための方法であって、各メタデータがデータファイルに対応している、少なくとも1つのメタデータのセットを処理することと、メタデータを分析するために要求されたパラメータ値を得ることと、得られたパラメータ値を用いてメタデータ要素のセット内の構造を判断することであって、複数のデータファイルの少なくともサブセットについて、パラメータ値を互いに用いてメタデータの少なくともサブセットを比較することによって構造が判断される、構造を判断することと、を含む、方法である。

40

【0027】

本発明の第12の態様は、第11の態様において、メタデータ要素の判断された構造を用いて、データファイルを群にクラスタ化することをさらに含む、方法である。

【0028】

50

本発明の第13の態様は、第12の態様において、データファイルの判断されたクラスタから境界を判断することをさらに含み、境界はデータファイルの判断されたクラスタ間に位置している、方法である。

【0029】

本発明の第14の態様は、第13の態様において、データファイルの1つのクラスタ内のメタデータ要素の少なくともいくつかを、データファイルの要素クラスタ内のメタデータ要素の少なくともいくつかと比較することによって、類似性値を判断することと、データファイルの1つのクラスタ内のメタデータ要素の少なくともいくつかを、データファイルの他のクラスタ内のメタデータ要素の少なくともいくつかと比較することによって、非類似性値を判断することと、をさらに含む、方法である。

10

【0030】

本発明の第15の態様は、第14の態様において、類似性値と非類似性値との相違に基づいて、データファイルのクラスタの要求された群に対応するパラメータ値を判断することをさらに含む、方法である。

【0031】

本発明の第16の態様は、データ処理装置上で実行可能であると共に、少なくとも各データファイルに関連した少なくとも1つのメタデータ要素を有するメタデータを用いることによって複数のデータファイルを構成するために用いられ得る、プログラムであって、プログラムが、データファイルの少なくともいくつかについて、そのデータファイルに関連した少なくとも1つのメタデータ要素を抽出するための命令と、抽出されたメタデータ要素についての値に基づいて、抽出されたメタデータ要素を要求された順序に構成するための命令と、パラメータ値を入力するための命令と、抽出されたメタデータ要素及び入力パラメータ値に基づいて、データファイルの少なくともいくつかを群に分割するための命令と、を含む、プログラムである。

20

【0032】

本発明の第17の態様は、第16の態様において、データファイルの少なくともいくつかを群に分割するための命令が、少なくとも1つのパラメータ値の少なくとも1つの各々について、抽出されたメタデータ要素の少なくともいくつか及びそのパラメータ値を用いて、複数のデータファイルの少なくとも2つについての類似性値を判断するための命令をさらに含む、プログラムである。

30

【0033】

本発明の第18の態様は、第17の態様において、少なくともいくつかのデータファイルの各々について、そのデータファイルについての及び多数の近くのデータファイルについての少なくとも1つの類似性値に基づいて、そのデータファイルについての少なくとも1つの新規性値を判断するための命令をさらに含む、プログラムである。

【0034】

本発明の第19の態様は、第17の態様において、少なくとも1つの類似性値を判断するための命令が、次の数式で少なくとも1つの類似性値を判断するための命令を含み、

$$S_K(i, j) = \exp\left(-\frac{|t_i - t_j|}{K}\right)$$

40

$S_K(i, j)$ が i^{th} データファイル及び j^{th} データファイルについての類似性値であり、 K がパラメータ値であり、 t_i 及び t_j が、 i^{th} 及び j^{th} データファイルについての少なくとも1つの抽出されたメタデータ要素の少なくとも一つのメタデータ要素の実値である、プログラムである。

【0035】

本発明の第20の態様は、第17の態様において、少なくとも1つの類似性値を判断するための命令が、次の数式で少なくとも1つの類似性値を判断するための命令を含み、

50

$$S_K(i, j) = \exp \left(\frac{1}{K} \left(\frac{\langle v_i, v_j \rangle}{\|v_i\| \|v_j\|} - 1 \right) \right)$$

$S_K(i, j)$ が i^{th} データファイル及び j^{th} データファイルについての類似性値であり、 K がパラメータ値であり、 v_i 及び v_j が、 i^{th} 及び j^{th} データファイルから判断された実際のベクトル値である、プログラムである。

【0036】

10

本発明の第21の態様は、第18の態様において、データファイルの少なくともいくつかについて判断された少なくとも1つの新規性値に基づいて、複数のデータファイルの境界位置間の少なくとも1つの境界位置を判断するための命令をさらに含む、プログラムである。

【0037】

本発明の第22の態様は、第18の態様において、少なくとも1つの新規性値を判断するための命令が、次の数式で少なくとも1つの新規性値を判断する命令を含み、

$$v_K(s) = \sum_{l,n=-5}^5 S_K(s+l, s+n) g(l, n)$$

20

$v_K(s)$ が新規性値であり、 g がガウシアンテーパーの 11×11 チェッカーボードカーネルである、プログラムである。

【0038】

本発明の第23の態様は、第21の態様において、判断された境界位置の少なくともいくつかについて、その境界位置についての信頼値を判断するための命令をさらに含む、プログラムである。

【0039】

本発明の第24の態様は、第23の態様において、少なくとも1つの信頼値を判断するための命令が、次の数式でこのような信頼値の各々を判断するための命令を含み、

30

$$C(B_K) = \sum_{l=1}^{|B_K|-1} \frac{1}{(b_{l+1} - b_l)^2} \sum_{i,j=b_l}^{b_{l+1}} S_K(i, j) - \sum_{l=1}^{|B_K|-2} \frac{1}{(b_{l+1} - b_l)(b_{l+2} - b_{l+1})} \sum_{i=b_l}^{b_{l+1}} \sum_{j=b_{l+1}}^{b_{l+2}} S_K(i, j)$$

$C(B_K)$ が B_K^{th} 境界についての信頼値であり、 $S_K(i, j)$ が i^{th} データファイル及び j^{th} データファイルについての類似性値であり、 b があるレベルの検出された境界である、プログラムである。

【0040】

本発明の第25の態様は、第23の態様において、本発明の判断された境界位置の少なくともいくつかについて、信頼値を最大にする少なくとも1つのパラメータ値の少なくとも1つを判断するための命令をさらに含む、プログラムである。

40

【0041】

本発明の第26の態様は、データファイルの対応するメタデータ要素と少なくとも関連した、少なくとも1つのメタデータ要素を有するメタデータを用いて、複数のデータファイルを構成するために用いられ得る、データファイル構成システムであって、データファイルの少なくともいくつかについて、データファイルと関連した少なくとも1つのメタデータ要素を抽出する、メタデータ抽出回路、ルーティン、又はアプリケーションと、抽出されたメタデータについての値に基づいて、抽出されたメタデータ要素を要求された順序に構成するための、メタデータ構成回路、ルーティン、又はアプリケーションと、少なくとも1つのパラメータ値の少なくとも1つについて、抽出されたメタデータ要素の少なく

50

ともいくつか及びそのパラメータ値を用いて、複数のデータファイルの少なくとも2つについての類似性値を判断する、類似性値判断回路、ルーティン、又はアプリケーションと、データファイルについての及び多数の近くのデータファイルについての少なくとも1つの類似性値に基づいて、そのデータファイルについての少なくとも1つの新規性値を判断する、新規性値判断回路、ルーティン、又はアプリケーションと、データファイルの少なくともいくつかについて判断された少なくとも1つの新規性値に基づいて複数のデータファイルの境界位置間の少なくとも1つの境界位置を判断することによって、抽出されたメタデータ要素及び入力パラメータ値に基づいて、データファイルの少なくともいくつかを群に分割する、データ分割回路、ルーティン、又はアプリケーションと、判断された境界位置の少なくともいくつかについて、その境界位置についての信頼値を判断し、判断された境界位置の少なくともいくつかについて、データ分割回路、ルーティン、又はアプリケーションが、信頼値を最大にする少なくとも1つのパラメータ値をさらに判断する、信頼値判断回路、ルーティン、又はアプリケーションと、を備える、データファイル構成システムである。

10

【発明を実施するための最良の形態】

【0042】

本発明のこれらの又は他の特徴及び利点は、本発明による方法及び装置の種々の例示的な実施の形態の次の詳細な説明に記載されており、又はこれらから明らかである。

【0043】

本発明の種々の例示的な実施の形態を、添付の図を参照して詳細に説明する。

20

【0044】

本発明によるシステム及び方法の種々の例示的な実施の形態の次の詳細な説明は、データファイルに対応するメタデータの処理に基づいて要求されたデータを構成することを主眼にしている。しかしながら、当然のことながら、本発明は開示された例示的な実施の形態のみに限定されない。一般には、本発明は、対応するメタデータを用いて多数のデータを構成する、任意の方法又は装置に用いられ得る。なお、本発明は、コンピュータを使用して実施されてもよい。また、そのようなコンピュータが端末として複数備わるネットワークシステム上で実施されてもよい。本発明の実施の形態におけるコンピュータは、少なくとも演算処理を行うプロセッサ、データ及びインストラクションをユーザが入力するための入力手段、データ及び処理結果を出力する出力手段、及び、データ及びプログラムを記憶する記憶手段を備える。以下に説明する本発明の実施の形態のシステムまたは方法は、当該システムまたは方法を実行するためのプログラムを記憶手段に記憶し、当該記憶手段から当該プログラムを読み出し、該プロセッサにより実行するようにしてもよい。

30

【0045】

図1は、本発明によるデータを構成するための方法の1つの例示的な実施の形態を概説する、フローチャートである。種々の例示的な実施の形態において、図1で概説された方法は、複数のデータファイル内の及び/又は関連したメタデータに基づいて、任意の要求されたデータの種類の複数のデータファイルを構成するために、用いられ得る。

【0046】

図1で示されるように、本方法の動作は、ステップS100で始まり、S200に続き、各データファイルのメタデータの少なくとも1つの要素が、構成されるべき複数のデータファイルから抽出される。次に、ステップS300では、抽出されたメタデータの要素は、抽出されたメタデータ要素についての1つ以上の値に基づいて、セットに構成され、たとえば、セット内での要求された順序及び識別という、指示が与えられる。動作は、次に、ステップS400に続く。

40

【0047】

ステップS400では、パラメータKについての値が選択される。次に、ステップS500では、メタデータが、要求通りに階層的に構成される。動作は、次に、ステップS600に続き、本方法の動作が終了する。

【0048】

50

当然のことながら、種々の例示的な実施の形態において、たとえば、メタデータの少なくとも1つの抽出された要素がタイムスタンプ要素を含む場合には、抽出されたメタデータ要素は時間的順序で構成されてもよい。あるいは、メタデータの少なくとも1つの抽出された要素がファイル名又は他のテキストストリングを含む場合には、メタデータ要素はアルファベット順で構成されてもよい。他の種々の例示的な実施の形態において、メタデータの少なくとも1つの抽出されたメタデータ要素が数値データを含む場合には、メタデータ要素は数値順に構成されていてもよい。さらに他の種々の例示的な実施の形態において、メタデータの少なくとも1つの抽出されたメタデータ要素は、たとえばGPSデータ等の位置を定義してもよい。当然のことながら、上述の、時間による、アルファベットによる、数値による、及び/又は位置によるメタデータ要素に加えて又はその代わりに、任意の他の適切なメタデータ要素が、構成する特徴として用いられ得る。これも当然のことであるが、選択されたメタデータ要素の値を順序付け又は構成する、任意の周知又は今後開発される方法が、データファイルを要求された順序に構成するために用いられてもよい。

10

20

30

40

50

【0049】

種々の例示的な実施の形態において、各抽出されたメタデータ要素は要求された識別を与えられるか、又はインデックスが付される。その結果、このような例示的な実施の形態では、各データファイルは、時間、名前、又は位置によって構成するメタデータ構成要素の実際の値に基づいてではなく、データファイルのセット内のそのメタデータ要素の値の位置によって、このように識別される。いいかえれば、たとえば、データファイルのセットは、タイムスタンプメタデータ要素の値に基づいて、時間的順序で構成される。しかしながら、データファイルは、次に、タイムスタンプメタデータ要素の絶対的時間値によってではなく、タイムスタンプメタデータ要素の時間値の観点から、データファイルのセット内に位置した順序によって、識別されるか、又はインデックスが付される。それにもかかわらず、各データファイルについてのメタデータ要素は、その絶対値を保持し続け、後に比較し得る。

【0050】

種々の例示的な実施の形態において、パラメータKは、数値を有する。パラメータKについての入力値は、デフォルト値又は要求値の場合がある。種々の例示的な実施の形態において、パラメータKは、セット内のデータファイルの各対の選択されたメタデータ要素又はセット内のデータファイル対のサブセット間の対の二つ一組の比較を行なうために、クラスタリング感度を判断する値である。それゆえ、パラメータKのより大きな値は、結果としてデータファイルのより粗いクラスタリングになる、比較を示す。言い換えれば、パラメータKのより大きな値は、分離したクラスタに分類されるために、メタデータについての値が互いにより離れていることを要求する。他方、パラメータKについてのより小さな値は、パラメータKについてのより大きな又はより小さな値のいずれかにおいて、多かれ少なかれ明らかになるメタデータの特有の特徴を統合するか又は強調するために、調整され得る。

【0051】

たとえば、パラメータKについてのより小さな値は、典型的には、非常に細かく離れた値、又はより小さな差異でより明らかになるメタデータの特徴を有するメタデータ要素について、より適切である。対照的に、パラメータKについてのより大きな値は、典型的には、非常に粗く離れた値又はより大きな差異でより明らかになるメタデータの特徴を有するメタデータ要素について、より適切である。その結果、パラメータKに対する要求値は、メタデータの種類と、メタデータの間隔と、セット内のメタデータ要素の数とに従って、異なるであろう。それゆえ、種々の例示的な実施の形態において、パラメータKについての複数の値は、メタデータを完全に分析して比較するために用いられる。このようにして、本発明による種々の例示的な実施の形態において、メタデータ要素の入力セットの先天的な分布に関して、仮定は行なわれない。パラメータKについてのこのような値を用いて分類及び/又は比較され得るメタデータの種々の例示的な種類は、たとえば、低レベル

画像特徴、GPSデータ、時間、月、及び/又は年におけるタイムスタンプを含む。

【0052】

図2は、ステップS500の要求されたメタデータを階層的に構成するための方法の1つの例示的な実施の形態をより詳細に概説する、フローチャートである。種々の例示的な実施の形態において、図2で概説された方法は、データファイルの任意の要求されたセットが、そのメタデータを用いることによって構成するために、用いられ得る。

【0053】

図2で示されたように、本方法の動作は、ステップS500で始まり、ステップS510に続き、ここでパラメータKについての値のリストが得られる。次に、ステップS520で、パラメータKについての値のリストから、最初の又は次の値が選択される。動作は、次に、ステップS530に続く。 10

【0054】

パラメータKについての値のリストは、ステップS400で選択されたパラメータKについての値に対応する。種々の例示的な実施の形態において、パラメータKについての複数の異なる値を含む、パラメータKについての値のリストは、たとえばメタデータ値のクイックスキャンに基づいて無作為に自動的に生じ得るか、又は手動で入力され得るかのいずれかである。種々の例示的な実施の形態において、リスト内のパラメータKについての値は、パラメータKについての複数の値を含む。

【0055】

ステップS530では、リスト内のインデックスが付されたメタデータ要素の各対について類似性値 S_K を得るために、リスト内のパラメータKについての値の各々が用いられる。 20

$$S_K(i, j) = \exp\left(-\frac{|t_i - t_j|}{K}\right) \quad (1)$$

ここで、 $S_K(i, j)$ は i^{th} 及び j^{th} データファイルについての類似性値であり、 K はパラメータKの値であり、 t_i 及び t_j は、 i^{th} 及び j^{th} データファイルの選択されたメタデータ要素の実際の値である。 30

【0056】

パラメータKについての特定の値を用いる、メタデータ要素の各比較された対についての類似性値 S_K の集合は、類似性行列として表現され得る。

【0057】

言い換えれば、 i^{th} 及び j^{th} データファイルのメタデータ要素の値 t_i 及び t_j についての類似性値 S_K を得るために、 i^{th} 及び j^{th} データファイルについてのメタデータは、パラメータKに基づいて比較し得る。 t 値がメタデータの実際の値であるから、1つの例示的な実施の形態において、メタデータがタイムスタンプである場合には、 t は分単位の時間であってもよい。 40

【0058】

類似性値 S_K を得るために用いられ得るメタデータ要素の実際の値の種類は、時間等のスカラー値である必要はない。類似性値 S_K を得るために、他の種類のメタデータ要素が用いられ得る。種々の例示的な実施の形態において、コンテンツベースの特徴ベクトルが、メタデータと共に、又は代わりに用いられ得る。この場合には、類似性値は次のとおりである。

$$S_K(i, j) = \exp \left(\frac{1}{K} \left(\frac{\langle v_i, v_j \rangle}{\|v_i\| \|v_j\|} - 1 \right) \right) \quad (2)$$

ここで、 v_i 及び v_j は、 i 及び j データファイルの選択されたメタデータ要素についての実際のベクトルである。他の適切な種類の値及び数式が、種々の他の例示的な実施の形態において用いられてもよい。動作は、次に、ステップ S 5 4 0 に続く。

【0059】

ステップ S 5 4 0 では、パラメータ K についての特定の値について生じた類似性行列 S_K の各要素について、新規性スコア v_K が得られる。新規性シェア v_K を得るための 1 つの方法は、類似性行列 $S_K(i, j)$ の主要な対角線 $S(i, j)$ に沿ってカーネルを相関させるために、適合したフィルタ技術を用いることである。すなわち、種々の例示的な実施の形態において、新規性スコア v_K は、類似性行列 S_K の対角線に沿ってのみ判断される。メタデータの群間の実際の境界を見出すために、種々の例示的な実施の形態において、新規性スコア $v_K(s)$ を計算するために、ガウシアンテーパの 11×11 チェッカーボードカーネル g が、次のように用いられる。

$$v_K(s) = \sum_{l,n=-5}^5 S_K(s+l, s+n) g(l, n) \quad (3)$$

ここで、 $v_K(s)$ は、パラメータ K についての特定の値についての類似性行列 S_K の i 要素及びガウシアンテーパの 11×11 チェッカーボードカーネル g についての新規性スコアである。

【0060】

数式 (3) では、 11×11 行列のゆえに、 -5 と $+5$ との間の l 及び n の範囲についての値が用いられる。種々の例示的な実施の形態において、たとえば、 9×9 行列等、他の寸法の行列が用いられてもよく、ここでは、 -4 と 4 との間の j 及び k の範囲についての値である。新規性スコア v_K を得るために、任意の要求された寸法のチェッカーボードカーネルを用いてもよい。

【0061】

チェッカーボードカーネルを用いることによって、完全な分析を実行する必要はない。むしろ、カーネルと同じ幅を有する主要な対角線の周囲のストリップのみを得る必要があり、それによって、データファイルの数に直線的に一致する、計算の複雑さを減らす。注目すべきことは、データのすべての可能な対ではなく、データの対のサブセットのみの比較が、任意の二つ組の比較において用いられてもよいことである。一般には、すべての可能な対のサブセットのみを用いることが、最小の性能劣化を伴う実質的な計算の節約につながる。

【0062】

パラメータ K の種々の値について新規性スコア v_K が判断される時には、新規性スコア内にいくつかのピークが現れる。注目すべきことは、パラメータ K の異なる値について異なるピークが現れることである。パラメータ K についての値が構造の範囲を示すから、パラメータ K についての異なる値は、類似性行列 S_K が異なる解像度での構造を明らかにすることを可能にする。新規性スコア v_K 内のピークは、次に、他の群と類似又はより近接したメタデータ要素値を有する、接続したデータ群間の境界の階層的なセット、すなわちクラスタを示す。それゆえ、新規性スコア v_K 内のピークは、類似したメタデータ値を有する群間の境界であり、他のクラスタから分離可能なメタデータ値のクラスタを示す。それゆえ、メタデータの群間の境界である、新規性スコア v_K 内のピークが得られる。動作は、次に、ステップ S 5 5 0 に続く。

10

20

30

40

50

【 0 0 6 3 】

ステップ S 5 5 0 では、最初に、パラメータ K の各値について新規性スコア v_k 内にすべてのピークを配置し、次に、検出された境界上に階層的構造を適用することによって、パラメータ K の各異なる値についての境界リストが得られる。種々の例示的な実施の形態において、パラメータ K の値のリスト内の各値を用いて、より粗いスケールからより細かいスケールまで、すなわちパラメータ K についての値を減少して、境界リストを得るための分析が行なわれる。ピーク値又は境界の階層的なセットを作るために、検出されたすべての境界を含むであろう境界リスト $B_k = [b_1, \dots, b_{n_k}]$ を用いて、パラメータ K の各値についての新規性スコア v_k 内のすべてのピークが、次に集められる。すなわち、粗いスケールすなわちパラメータ K のより大きい値で検出されたすべての境界は、すべてのより細かいスケールすなわちパラメータ K のより少ない値についての境界リスト内に含まれるであろう。より粗いスケールで得られたより遠く離れた群間の境界は、依然としてより細かいスケールで存在すると仮定される。

【 0 0 6 4 】

新規性スコア v_k が、局所的な最大値にあって、類似性測定の最大及び類似性行列の主要な対角線に沿って相関されたカーネルから判断されるところに、境界が位置する。新規性スコアの最大又は最小を得る他の方法は、たとえば数式 (3) の導関数を得ることである。動作は、次に、ステップ S 5 6 0 に続く。

【 0 0 6 5 】

ステップ S 5 6 0 では、パラメータ K の各値について、類似性値 S_k と、新規性スコア v_k と、境界 b_k とを得ることによって、境界を判断するためにリスト内のパラメータ K についてのすべての値が用いられたかが判断される。否の場合には、動作はステップ S 5 2 0 に戻る。さもなければ、動作は、ステップ S 5 7 0 に続く。

【 0 0 6 6 】

ステップ S 5 7 0 では、境界 B_k のリストによって表された検出された境界は、検出された境界の階層内の各レベルについてランクされて来たクラスタリングの結果を表わす、信頼スコア $C(B_k)$ を得るために用いられる。信頼スコア $C(B_k)$ は、次に数式によって示された平均的なクラス内の類似性とクラス間の非類似性とに基づく。

$$C(B_k) = \sum_{l=1}^{|B_k|-1} \frac{1}{(b_{l+1} - b_l)^2} \sum_{i,j=b_l}^{b_{l+1}} S_k(i, j) - \sum_{l=1}^{|B_k|-2} \frac{1}{(b_{l+1} - b_l)(b_{l+2} - b_{l+1})} \sum_{i=b_l}^{b_{l+1}} \sum_{j=b_{l+1}}^{b_{l+2}} S_k(i, j) \quad (4)$$

ここで、 $C(B_k)$ は信頼スコアであり、

b は各レベルにおける検出された境界である。

【 0 0 6 7 】

上述のように、各クラスタ内のデータファイル間の平均的なクラス内の類似性を定量化する第 1 の合計と、隣接したクラスタ内のデータファイル間の平均的なクラス間の類似性を定量化する第 2 の合計とは、クラスタ間の非類似性を定量化するために、無効にされる。第 1 の合計及び第 2 の合計についての変化率は、パラメータ K の値に従って変化する。それゆえ、パラメータ K についての複数の値については、1 つの値は、信頼スコア $C(B_k)$ が最大化されることを可能にするであろう。その結果、動作はステップ S 5 8 0 に続き、そこでは、信頼スコア $C(B_k)$ を最大にするパラメータ K の値についての境界リスト B_k が得られる。次に、動作は、ステップ S 5 9 0 に続き、ここで動作がステップ S 6 0 0 に戻る。ベイズ情報基準 (BIC) 等の信頼スコア $C(B_k)$ を得るために、他の種類の統計的測定が用いられ得る。ベイズ情報基準のいくつかの例が、D. ハッカーマン (D. Heckermann) 著「ベイズのネットワークを学ぶに当たっての指導書」(“A tutorial on learning with Bayesian networks”), 技術レポート (Technical Report) M S R - T R - 9 5 - 0 6, マイクロソフト・リサーチ (Microsoft Research)、レッドモンド、ワシントン (1995 年。1996 年改定) (非特許文献 1)、S. チェンら (S. Chen et al.) 「ベイズの情報基準による、話者、環境、及びチャネル変化の検出及びクラスタ

リング」(“Speaker, environmental and change detection and clustering via the Bayesian information criterion”) D A R P A 音声認識ワークショップ (DARPA Speech Recognition Workshop) (1998年) (非特許文献6)、及びS. レナルら (S. Renal et al.) 著「会議室からのオーディオ情報アクセス」(“Audio Information Access from Meeting Rooms”) (2003年4月) (非特許文献7) で述べられており、それらの各々は、本明細書にその全体を参照することによって組み込まれる。

【0068】

本発明によるシステム及び方法の1つの例示的な使用は、階層的なクラスタリングによってデジタル写真をタイムベースのイベントに構成することを含む。デジタルカメラの急増と共に、パーソナルコンピュータ上に蓄積するデジタル写真の数が急速に増えている。典型的にはJ P E G画像ファイルフォーマットである個々のデジタル画像ファイルは、典型的にはExif (エグジフ、Exchangeable Image File Format) で保存されている、デジタルファイル内の大量のメタデータを含む。このようなメタデータは、いつ写真が撮影されたか又はその後いつ再保存若しくは修正されたかを示すタイムスタンプを含む。それにもかかわらず、複数のメタデータは画像ファイルと共に記録されてもよいから、オリジナルのタイムスタンプ、又は任意のその後に修正されたタイムスタンプ等の情報は、メタデータとして別々に記録されてもよく、本発明によるシステム及び方法の種々の例示的な実施の形態を用いて、個別に抽出及び分析可能である。

10

【0069】

1つの例示的な実施の形態において、512の写真のクラスタリングが用いられた。最初に、すべての写真がタイムスタンプ(メタデータ)を有していて、撮影者によって、意味のあるフォルダ、すなわち特有のイベント、の中に手動で置かれた。これらの写真のこの手動のクラスタリングは、グランドトルースクラスタリング(ground truth clustering)として、次の説明において参照されるであろう。

20

【0070】

各写真についてのExifヘッダが、その写真についてのタイムスタンプを抽出するために、最初に処理された。抽出されたタイムスタンプは、最初に構成され、時間で順序付けられた。タイムスタンプは、分(minutes)等の任意の基本的な時間単位を用いて、時間順に順序付けられた。しかしながら、一旦タイムスタンプが時間順に順序付けられると、次に、各タイムスタンプ及びこのような各対応する写真は、インデックス又は時間順の番号若しくは値を付与され、引き続いてその後、タイムスタンプの絶対的な時間値によってではなく、このインデックスによって参照された。

30

【0071】

タイムスタンプを抽出して写真を構成するための最初の処理の後に、タイムスタンプの集合の構造は、類似性行列 S_k を作ることによって評価された。図3は、グランドトルースクラスタリングから生じた類似性行列 S_k について得られた結果を図示的に示す。図3で図示的に表現された類似性行列 S_k の要素についての値は、同一のフォルダからの写真の対については1であり、撮影者によって異なるフォルダ内に保存された写真の対については0である。写真は、上述のように、時間順にインデックスが付される。類似性行列 S_k の (i, j) 要素についての値を判断するために、その中に i^{th} 及び j^{th} 写真が保存されたフォルダの名前が比較される。それらが同じ場合には、 (i, j) 要素は、1の値が割り当てられる。さもなければ、それは、0の値が割り当てられる。種々の例示的な実施の形態において、行列の主要な対角線に沿った類似性行列 S_k の要素のブロックは、各フォルダ内の写真の群と対応する。

40

【0072】

図3で示された類似性行列 S_k の主要な対角線に沿ったチェッカーボードパターンは、すでに異なるイベントに分類された写真を含むフォルダ間の境界を示す。それゆえ、チェッカーボードパターンは、異なるイベントの写真の群間の時間的順序での境界の図示的な表現である。類似性行列の i^{th} 及び j^{th} 要素として写真が表現されたときには、チェッカーボードパターンは、写真が記述するイベントも時間においてまとまりがない一方、類似

50

性行列内で写真が隣接していることを示す。

【0073】

図4は、グラントルースクラスタリングについて生じた新規性スコア v_k を示す。新規性スコア v_k は、ガウシアンテーパーの 11×11 チェッカーボードカーネル g を用いて得られる。図4は、図3で示されたチェッカーボードに対応する新規性スコア v_k のピークを示す。たとえば、図3では、2つの黒い正方形によって示された2つの比較的大きな群は、インデックス値210の近くで分離される。2つの正方形は、インデックス値210の近くで接触するだけである。2つの正方形が単に接触する位置では、写真の2つの群間の境界を示す。図4では、この境界を示すインデックス値210の近くの新規性スコア v_k に、対応するピークがある。

10

【0074】

図5～10は、グラントルースクラスタリングでクラスタ化された写真を用いて、 10^3 分、 10^4 分、 10^5 分のパラメータ K の値について得られた、いくつかの類似性行列 S_k 及びそれらの対応する新規性スコア v_k を示す。図5、7、9は、それぞれ、 10^3 分、 10^4 分、 10^5 分のパラメータ K の値についての類似性行列 S_k を示す。図6、8、10は、それぞれ、 10^3 分、 10^4 分、 10^5 分のパラメータ K の値についての新規性スコア v_k を示す。パラメータ K についての3つの異なる値は、3つの異なる解像度を表す。特に、パラメータ K についての値が少ない程、解像度は大きくなり、ここではタイムスタンプの群間のより細かい非類似性が明らかになる。

【0075】

20

図5、7、9で示されたように、類似性行列 S_k は、異なる解像度での構造を明らかにする。それにもかかわらず、パラメータ K についてのより大きな値では、詳細は、パラメータ K についてのより小さい値について程容易には現れない。パラメータ K についての値を用いている極端な例が、図12及び13で示される。2つの異なる類似性行列で現れる例示的な写真のインデックスを用いて、図12は、パラメータ K ($K=10$) についての10の値について得られた類似性行列の部位を示す。図13は、パラメータ K ($K=1,000$) についての1,000の値について得られた類似性行列の部位を示す。図12及び13で示されたように、パラメータ K についてのより大きな値について得られるよりも、パラメータ K についてのより小さな値については、より良い境界の定義が得られる。これは、数式(1)によって、境界のいずれかの側におけるクラスタ内の写真は、パラメータ K の異なる値については異なるクラス内の類似性を示すことから生じる。これは、次に、チェッカーボードカーネルとの相関の強度を変化させる。それゆえ、類似性測定 S_k は、低レベル画像特徴、GPSデータ、又は他のメタデータ等の他の特徴を統合又は強調するために、加除可能である。

30

【0076】

上記のように、異なる特徴は、パラメータ K の異なる値で、より明らかになる。対応する新規性スコア v_k では、境界点は、分析のスケール、すなわちパラメータ K の値に応じて著しく変化する。図6、8、10では、パラメータ K の値の限られた数についての新規性スコア v_k が示されている。しかしながら、図11では、パラメータ K の値のより大きな数についての新規性スコア v_k が示されている。図11で示されたように、新規性スコア v_k は、パラメータ K の値と共に大きく変化し、新規性スコア v_k は、異なるスケールすなわちパラメータ K の値で、異なる境界ピークを示す。これは、異なるイベントが異なる時間範囲を有するから生じる。すなわち、休暇又は誕生パーティ等のイベントは、異なる時間範囲を有するであろう。たとえば、後者のイベントは、一般には前者のイベントの時間的範囲に比べてより短い時間範囲を有するであろう。

40

【0077】

図11では、最小の新規性スコア v_k は、 $S_{(k)}$ 内の高い自己類似性の、すなわち低い新規性の領域に対応する。このようにして、領域は、このような高い自己類似性の領域間に、優先的に位置している。境界は、パラメータ K の値を減らすことによって順序付けられ、検出された境界上に、階層的構造が与えられる。このような階層は、検出された境界上

50

に適用されてもよい。言い換えれば、より細かいスケールについての境界のセット内に、非常に粗いスケール（高い K 値）からのすべての検出された境界が含まれているところでは、階層的境界のセットが作られても良い。この技術を用いると、より重要でない境界がさらに検出されるにつれて、より重要な境界が保持されることを可能にする。

【0078】

本技術は、あるスケールでは、すなわちパラメータ K のある値については、検出されたイベント境界が最大の新規性スコアに近づくはずである、という仮定をもとにしている。パラメータ K の各値については、境界を示す新規性スコア v_K 内のピークは、第1の相違の分析によって検出される。所与のしきい値スコアを用いると、たとえば、同じイベントである写真中の時間値内にある異常に長いギャップのゆえに、現れるかも知れない擬制のピークを検出することを回避できる。このような所与のしきい値スコアは、最小しきい値スコアとして用いられてもよい。たとえば、5よりも大きい新規性スコアが、各隣接範囲におけるピークとして選択され得る。

10

【0079】

図14は、数式(4)によって表されたように、各クラスタ内の選択されたメタデータ要素についての値間の平均的なクラス内の類似性と、隣接したクラスタ内の選択されたメタデータ要素についての値間の平均的なクラス間の類似性との相違である、推測されたクラスタ内の信頼を定量化する概念を示す。クラス内の類似性条件は、主要な対角線に沿った領域の条件全体の平均である。クラス間の類似性条件は、主要な対角線から離れた矩形領域の平均である。図14は、信頼スコアの計算を図で示す。

20

【0080】

この信頼測定 $C(B_K)$ は、検出されたクラスタの数とパラメータ K の値との両方に明らかに依存する。図15~17は、動作を表す。図15~17は、数式(4)で定義された信頼測定を形成するために平均化され、合計されたそれぞれの類似性行列 S_K の領域を示す。図15は、パラメータ K ($K = 1778.28$) についての1778.28の値についての行列を示す。図16は、 $K = 1, 000$ についての行列を示す。最後に、図17は、 $K = 562.34$ についての行列を示す。図15~17で示された行列の表現では、 $C(B_K)$ に貢献しない要素は、行列内では0にセットされている。図15~17では、パラメータ K のより大きな値については、パラメータ K についてのより低い値についても、より低い信頼スコアが得られる。たとえば、 $K = 1, 000$ (図16) については、信頼スコア $C(B_K)$ は21.09886であり、これは $K = 1778.28$ (図15) についての11.7814の信頼スコア $C(B_K)$ よりも大きい。実際、図16は、より少ない数のクラスタであって、比較的低い類似性についてのクラスタ化された領域を示す。他方、図17の $K = 562.34$ についての行列は、図16の $K = 1, 000$ についての行列よりも多いクラスタを示す。しかし、パラメータ K の値がより小さいゆえに、低い類似性の領域がクラスタ化される。このようにして、当然のことながら、種々の例示的な実施の形態において、信頼測定によって、類似性分析についての1つの適切なスケールが強調される。

30

【0081】

図18は、本発明によるデータ構成システム100の1つの例示的な実施の形態の構成図である。図18で示されたように、データ構成システム100は、1つ以上の制御及び/又はデータバス及び/又はアプリケーションプログラミングインタフェース195によって相互接続された、入力/出力インタフェース110、コントローラ120、メモリ130、メタデータ抽出回路、ルーティン、又はアプリケーション140、メタデータ構成回路、ルーティン、又はアプリケーション150、類似性値判断回路、ルーティン、又はアプリケーション160、新規性値判断回路、ルーティン、又はアプリケーション170、データ分割回路、ルーティン、又はアプリケーション180、及び信頼値判断回路、ルーティン、又はアプリケーション190を含む。

40

【0082】

図18で示されたように、表示装置102、1つ以上のユーザ入力装置106、データ

50

送信装置 200、及びデータ受信装置 220 が、リンク 104、108、210、230 によって、それぞれデータ構成システム 100 に接続されている。

【0083】

一般には、図 18 で示されたデータ送信装置 200 は、データファイル及びそれらの対応するメタデータをデータ構成システム 100 に供給することができる、任意の周知又は今後開発される装置であってもよい。一般には、図 18 で示されたデータ受信装置 220 は、データ構成システム 100 からの任意のデータを受信できる、任意の周知又は今後開発される装置であってもよい。

【0084】

データ送信装置 200 及び / 又はデータ受信装置 220 は、データ構成システム 100 と一体化されていてもよい。さらに、データ構成システム 100 は、捉えられた写真を自動的にフォルダ内に構成するデジタルカメラ等、多機能を実行するより大きなシステム内で、データ送信装置 200 及び / 又はデータ受信装置 220 に加えて追加の機能を与える装置と一体化されていてもよい。

【0085】

それぞれの 1 つ以上のユーザ入力装置 106 の各々は、キーボード、マウス、ジョイスティック、トラックボール、タッチパッド、タッチスクリーン、ペンベースのシステム、マイクロホン、及び関連した音声認識ソフトウェア、又はデータ構成システム 100 にデータ及び / 又はユーザコマンドを入力するための任意の他の周知又は今後開発される装置等の、複数の入力装置の 1 つ又は任意の組合せであってもよい。当然のことながら、図 18 の 1 つ以上のユーザ入力装置 106 は、同じ種類の装置である必要はない。

【0086】

表示装置 102、1 つ以上のユーザ入力装置 106、データ送信装置 200、及びデータ受信装置 220 をデータ構成システム 100 に接続しているリンク 104、108、210、230 のリンクの各々は、信号線、直接ケーブル接続、モデム、ローカルエリアネットワーク、広域ネットワーク、イントラネット、インターネット、任意の他の分散型の処理ネットワーク、又は任意の他の周知又は今後開発される接続装置又は構造であってもよい。当然のことながら、これらのリンク 104、108、210、230 のいずれも、有線又は無線の部位を含んでもよい。一般には、リンク 104、108、210、230 の各々は、それぞれの装置をデータ構成システム 100 に接続するために使用可能な、任意の周知又は今後開発される接続システム又は構造を用いて実行され得る。当然のことながら、リンク 104、108、210、230 は、同じ種類である必要はない。

【0087】

図 18 で示されたように、メモリ 130 は、可変、揮発、若しくは不揮発メモリ又は非可変若しくは固定したメモリの、任意の適切な組合せを用いて実行し得る。可変メモリは、揮発又は不揮発のいずれでも、任意の 1 つ以上の、スタティック又はダイナミック RAM、フロッピディスク及びディスクドライブ、書込可能又は再書込可能な光ディスク及びディスクドライブ、ハードディスクドライブ、フラッシュメモリ等を用いて、実行され得る。同様に、非可変又は固定のメモリは、CD ROM 若しくは DVD-ROM ディスク及びディスクドライブ等、任意の 1 つ以上の ROM、PROM、EPROM、EEPROM 及び光 ROM ディスクを用いて、実行され得る。

【0088】

データ構成システム 100 の種々の実施の形態が、プログラムされた汎用コンピュータ、専用コンピュータ、マイクロコンピュータ等の上で実行するソフトウェアとして実行され得る。これも当然のことながら、図 18 で示された回路、ルーティン、及び / 又はアプリケーションの各々は、適切にプログラムされた汎用データプロセッサの一部位として実行され得る。あるいは、図 18 で示された回路、ルーティン、及び / 又はアプリケーションの各々は、ASIC、デジタルシグナルプロセッサ (DSP)、FPGA、PLD、PLA、及び / 又は PAL、又は個別論理素子又は個別回路素子内の物理的に異なるハードウェア回路として、実行し得る。一般には、同様に図 1 及び 2 で示されたフローチャート

10

20

30

40

50

を実行し得る有限状態の機械を実行し得る任意の装置が、データ構成システム 100 を実行するために用いられ得る。図 18 で示された、特定の形式の回路、ルーティン、アプリケーション、対象、及び / 又は管理者は、設計上の選択として解釈し、当業者には明らかでありかつ予期し得るものである。当然のことながら、図 18 で示された、回路、ルーティン、アプリケーション、対象、及び / 又は管理者は、同じ設計である必要はない。

【0089】

メタデータ抽出回路、ルーティン、又はアプリケーション 140 は、データファイルと関連した少なくとも 1 つのメタデータ要素を抽出する。各データファイルのメタデータの少なくとも 1 つの要素は、構成されるべき複数のデータファイルから抽出される。典型的には J P E G 画像ファイルフォーマットである、デジタル画像ファイル等のデータファイルは、典型的には標準的な交換可能な画像ファイルフォーマット (E x i f) で保存された、デジタルファイル内の大量のメタデータを含む。このような抽出可能なメタデータは、写真が撮影された時又は次に再保存若しくは修正された時を示す。

10

【0090】

メタデータ構成回路、ルーティン、又はアプリケーション 150 は、抽出されたメタデータ要素についての値に基づいて、抽出されたメタデータ要素を、要求された順序に構成する。抽出されたメタデータ要素は、時間の、アルファベット順の、数字の、及び / 又は位置の特徴等の、任意の要求された構成する特徴を用いて構成され、割り当てられた又はインデックスが付された識別値に基づいて、抽出されたメタデータ要素を順序付けることができる。

20

【0091】

類似性値判断回路、ルーティン、又はアプリケーション 160 は、少なくとも 1 つのパラメータ値の少なくとも 1 つについて、抽出されたメタデータ要素の少なくともいくつか及びそのパラメータ値を用いて、複数のデータファイルの少なくとも 2 つについての類似性値を判断する。それゆえ、類似性値判断回路、ルーティン、又はアプリケーション 160 は、データファイルの各々のこのような対の類似性値を得るために、パラメータ値を用いて、少なくとも 1 対のデータファイルについてのメタデータを比較する。

【0092】

新規性値判断回路、ルーティン、又はアプリケーション 170 は、複数の類似性値に基づいて、データファイルについての少なくとも 1 つの新規性値を判断する。すなわち、新規性値判断回路、ルーティン、又はアプリケーション 170 は、要求されたデータファイルの数についての類似性値に基づいて、新規性値を判断する。

30

【0093】

データ分割回路、ルーティン、又はアプリケーション 180 は、抽出されたメタデータ要素及び入力パラメータ値に基づいて、データファイルの少なくともいくつかを群に分割する。種々の例示的な実施の形態において、データ分割回路、ルーティン、又はアプリケーション 180 は、データファイルの少なくともいくつかについて判断された少なくとも 1 つの新規性値に基づいて複数のデータファイルの境界位置間の少なくとも 1 つの境界位置を判断することと、判断された境界位置の少なくともいくつかについて、信頼値を最大にする少なくとも 1 つのパラメータを判断することとによって、抽出されたメタデータ要素及び入力パラメータ値に基づいて、データファイルの少なくともいくつかを群に分割する。

40

【0094】

信頼値判断回路、ルーティン、又はアプリケーション 190 は、判断された境界位置の少なくともいくつかについて、その境界位置についての信頼値を判断する。

【0095】

動作中には、データ構成システム 100 は、それぞれが対応するメタデータを有する複数のデータファイルを入力するか又はさもなければ得て、入力パラメータについての値をリンク 210 を通じてデータ送信装置 200 から入力してもよく、及び / 又はメモリ 130 から 1 つ以上のデータファイルを読み出してもよい。入力パラメータは、ユーザ入力装

50

置 1 0 6 を通じて入力されてもよい。データ送信装置 2 0 0 から得られた場合には、入力 / 出力インタフェース 1 1 0 は、データファイル及び / 又は入力パラメータを入力し、コントローラ 1 2 0 の制御下で、任意の適切なデータファイルを、メタデータ抽出回路、ルーティン、又はアプリケーション 1 4 0 に送る。

【 0 0 9 6 】

メタデータ抽出回路、ルーティン、又はアプリケーション 1 4 0 は、入力データファイルの少なくともいくつかと関連した少なくとも 1 つのメタデータ要素を抽出する。メタデータ抽出回路、ルーティン、又はアプリケーション 1 4 0 は、次に、コントローラ 1 2 0 の制御下で、抽出されたメタデータ要素をメモリ 1 3 0 に保存するか、又は抽出されたメタデータ要素をメタデータ構成回路、ルーティン、又はアプリケーション 1 5 0 に直接出力する。メタデータ構成回路、ルーティン、又はアプリケーション 1 5 0 は、コントローラ 1 2 0 の制御下で、抽出されたメタデータ要素を入力し、抽出されたメタデータ要素についての値に基づいて、抽出されたメタデータ要素を要求された順序に構成する。メタデータ構成回路、ルーティン、又はアプリケーション 1 5 0 は、次に、コントローラ 1 2 0 の制御下で、順序付けられた抽出されたメタデータをメモリ 1 3 0 に保存するか、又は順序付けられた抽出されたメタデータ要素を類似性値判断回路、ルーティン、又はアプリケーション 1 6 0 に直接出力する。

10

【 0 0 9 7 】

類似性値判断回路、ルーティン、又はアプリケーション 1 6 0 は、コントローラ 1 2 0 の制御下で、順序付けられたメタデータ要素及び / 又は対応するデータファイルを入力し、少なくとも 1 つのパラメータ値の少なくとも 1 つについて、抽出されたメタデータ要素の少なくともいくつか及び又はそれらのデータファイルの内容及びそのパラメータ値を用いて、複数のデータファイルの 2 つで構成された対の少なくとも 1 対についての類似性値を判断する。類似性値判断回路、ルーティン、又はアプリケーション 1 6 0 は、次に、コントローラ 1 2 0 の制御下で、判断された類似性値をメモリ 1 3 0 に保存するか、又は、類似性値を、新規性値判断回路、ルーティン、又はアプリケーション 1 7 0 に直接出力する。

20

【 0 0 9 8 】

新規性値判断回路、ルーティン、又はアプリケーション 1 7 0 は、コントローラ 1 2 0 の制御下で、類似性値の少なくともいくつかを入力して、入力類似性値と関連した多数のデータファイルのそれぞれについて、そのデータファイルについての類似性値及び要求された数の周囲のデータファイルに基づいて、このようなデータファイルの各々についての少なくとも 1 つの新規性値を判断する。新規性値判断回路、ルーティン、又はアプリケーション 1 7 0 は、次に、コントローラ 1 2 0 の制御下で、判断された新規性値をメモリ 1 3 0 に保存するか、又は、判断された新規性値をデータ分割回路、ルーティン、又はアプリケーション 1 8 0 に直接出力する。

30

【 0 0 9 9 】

データ分割回路、ルーティン、又はアプリケーション 1 8 0 は、コントローラ 1 2 0 の制御下で、新規性値の少なくともいくつかを入力し、複数のデータファイルの種々の境界位置間の少なくとも 1 つの境界位置を判断することによって、データファイルの少なくともいくつかについて判断された少なくとも 1 つの新規性値に基づいて、対応するデータファイルを群に分割する。データ分割回路、ルーティン、又はアプリケーション 1 8 0 は、次に、コントローラ 1 2 0 の制御下で、判断された境界位置をメモリ 1 3 0 に保存するか、又は判断された境界位置を信頼値判断回路、ルーティン、又はアプリケーション 1 9 0 に出力する。

40

【 0 1 0 0 】

信頼値判断回路、ルーティン、又はアプリケーション 1 9 0 は、コントローラ 1 2 0 の制御下で、1 つ以上の境界位置を入力し、判断された境界位置の少なくともいくつかについて、判断された境界位置の少なくともいくつかについての境界位置についての信頼値を判断する。信頼値判断回路、ルーティン、又はアプリケーション 1 9 0 は、次に、コント

50

ローラ 120 の制御下で、判断された信頼値をメモリに保存するか、又は判断された信頼値をデータ分割回路、ルーティン、又はアプリケーション 180 に出力する。データ分割回路、ルーティン、又はアプリケーション 180 は、次に、判断された境界位置の少なくともいくつかについての信頼値を最大にする少なくとも 1 つのパラメータ値を判断する。それゆえ、データ構成システム 100 の動作中は、入力パラメータ値、抽出された順序付けられたメタデータ要素、及び / 又は対応するデータファイルの内容は、順序付けられた抽出されたメタデータ要素、及び / 又はデータファイルの対応する内容、及び入力パラメータ値に基づいて、読み出された / 受け取られたデータファイルの少なくともいくつかを用いて、群に構成される。分割され、このようにして構成されたデータファイルは、次に、さらにメモリ 130 内に保存され、データ受信装置 220 に出力され、及び / 又は表示装置 102 上に表示され得る。 10

【0101】

図 18 は、データ構成ユニット 100 を表示装置 102 から分離された装置として示しているが、ユーザ入力装置 106、データ送信装置 200、及び / 又はデータ受信装置 220、及びデータ構成システム 100 は、一体化された装置であってもよい。一体化された構成では、表示装置 102、ユーザ入力装置 106、データ送信装置 200、及び / 又はデータ受信装置 220 からの 2 つ以上のデータ構成システム 100 が、単一の装置に含まれてもよい。

【0102】

あるいは、データ構成システム 100 は、メタデータ抽出回路、ルーティン、又はアプリケーション 140、メタデータ構成回路、ルーティン、又はアプリケーション 150、類似性値判断回路、ルーティン、又はアプリケーション 160、新規性値判断回路、ルーティン、又はアプリケーション 170、データ分割回路、ルーティン、又はアプリケーション 180、及び信頼値判断回路、ルーティン、又はアプリケーション 190、コントローラ 120、メモリ 130、及び / 又は入力 / 出力インタフェース 110 を含む、分離された装置であってもよい。さらに、分離された回路、ルーティン、及び / 又はアプリケーションとして示されているが、メタデータ抽出回路、ルーティン、又はアプリケーション 140、メタデータ構成回路、ルーティン、又はアプリケーション 150、類似性値判断回路、ルーティン、又はアプリケーション 160、新規性値判断回路、ルーティン、又はアプリケーション 170、データ分割回路、ルーティン、又はアプリケーション 180、及び信頼値判断回路、ルーティン、又はアプリケーション 190 は、それ自身が、種々の組み合わせで一体化されていてもよい。 20 30

【0103】

本発明を上述の概説された例示的な実施の形態に関連して述べてきたが、周知又は現在では予見できないかも知れないにしろ、種々の代替物、修正、変更、改良、及び / 又は実質的な等価物が、少なくとも本技術分野における通常の知識を有する者には自明になるかも知れない。したがって、上述した本発明の例示的な実施の形態は、例示的であり、制限しないことを意図されている。本発明の精神及び範囲から逸脱せずに種々の変更が可能である。それゆえ、提出されて、それらが補正されるかも知れない特許請求の範囲は、すべての周知又は今後開発される代替物、修正、変更、改良、及び / 又は実質的な等価物を含むように意図されている。 40

【図面の簡単な説明】

【0104】

【図 1】本発明によるデータを構成するための方法の 1 つの例示的な実施の形態を概説するフローチャートである。

【図 2】本発明による要求されたデータを構成するための方法の 1 つの例示的な実施の形態をより詳細に概説するフローチャートである。

【図 3】類似性行列及び新規性スコアについて得られた結果の 1 つの例示的な実施の形態を図示する。

【図 4】類似性行列及び新規性スコアについて得られた結果の 1 つの例示的な実施の形態 50

を図示する。

【図 5】複数の類似性行列について得られた結果及びそれらの対応する新規性スコアの例示的な実施の形態を図示する。

【図 6】複数の類似性行列について得られた結果及びそれらの対応する新規性スコアの例示的な実施の形態を図示する。

【図 7】複数の類似性行列について得られた結果及びそれらの対応する新規性スコアの例示的な実施の形態を図示する。

【図 8】複数の類似性行列について得られた結果及びそれらの対応する新規性スコアの例示的な実施の形態を図示する。

【図 9】複数の類似性行列について得られた結果及びそれらの対応する新規性スコアの例示的な実施の形態を図示する。 10

【図 10】複数の類似性行列について得られた結果及びそれらの対応する新規性スコアの例示的な実施の形態を図示する。

【図 11】パラメータ K 値に応じて変化する境界について判断された新規性スコアの 1 つの例示的な実施の形態を図示する。

【図 12】2 つの異なるパラメータ K 値について判断された類似性行列の例示的な実施の形態を図示する。

【図 13】2 つの異なるパラメータ K 値について判断された類似性行列の例示的な実施の形態を図示する。

【図 14】信頼スコアの 1 つの例示的な実施の形態を示す。 20

【図 15】3 つの異なるパラメータ K 値についての類似性行列の例示的な実施の形態を図示する。

【図 16】3 つの異なるパラメータ K 値についての類似性行列の例示的な実施の形態を図示する。

【図 17】3 つの異なるパラメータ K 値についての類似性行列の例示的な実施の形態を図示する。

【図 18】本発明によるデータ構成システムの 1 つの例示的な実施の構成図である。

【符号の説明】

【0105】

S200：メタデータを抽出する 30

S300：抽出されたデータを順序付けられたセットに構成する

S400：パラメータ K についての値を判断する

S500：要求されたデータを構成する

100：データ構成システム

140：メタデータ抽出回路、ルーティン、又はアプリケーション

150：メタデータ構成回路、ルーティン、又はアプリケーション

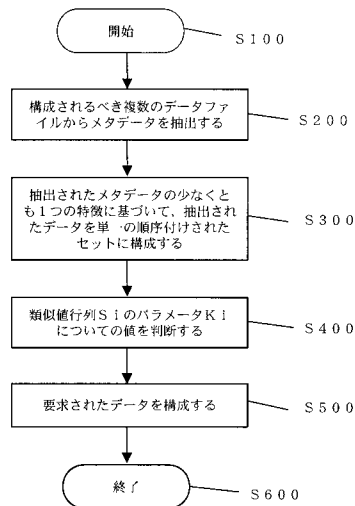
160：類似性値判断回路、ルーティン、又はアプリケーション

170：新規性値判断回路、ルーティン、又はアプリケーション

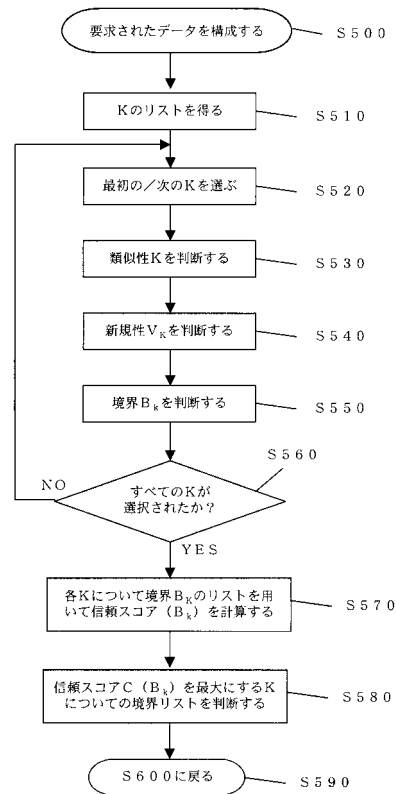
180：データ分割回路、ルーティン、又はアプリケーション

190：信頼値判断回路、ルーティン、又はアプリケーション 40

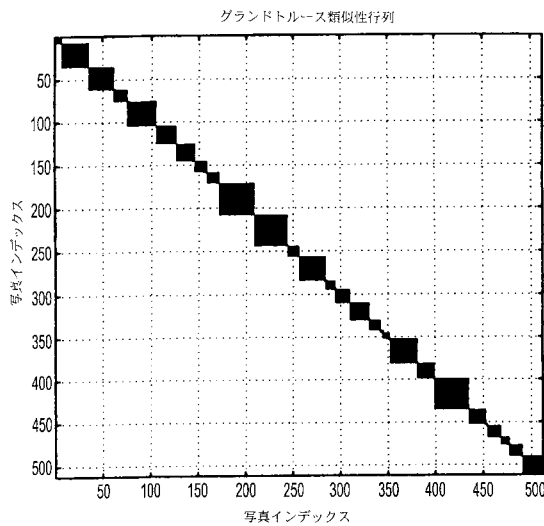
【図 1】



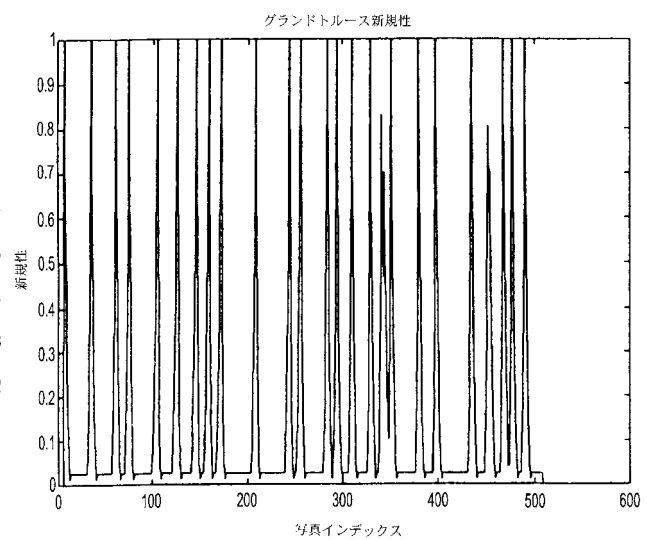
【図 2】



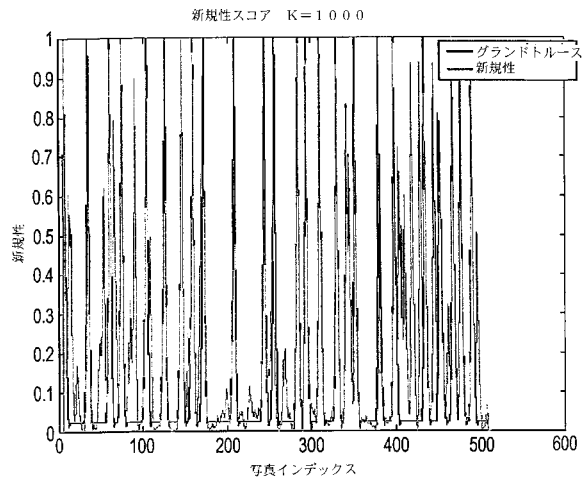
【図 3】



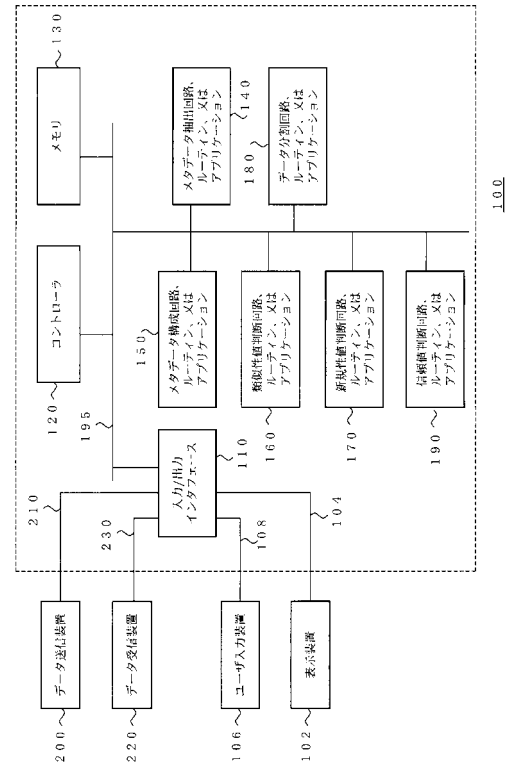
【図 4】



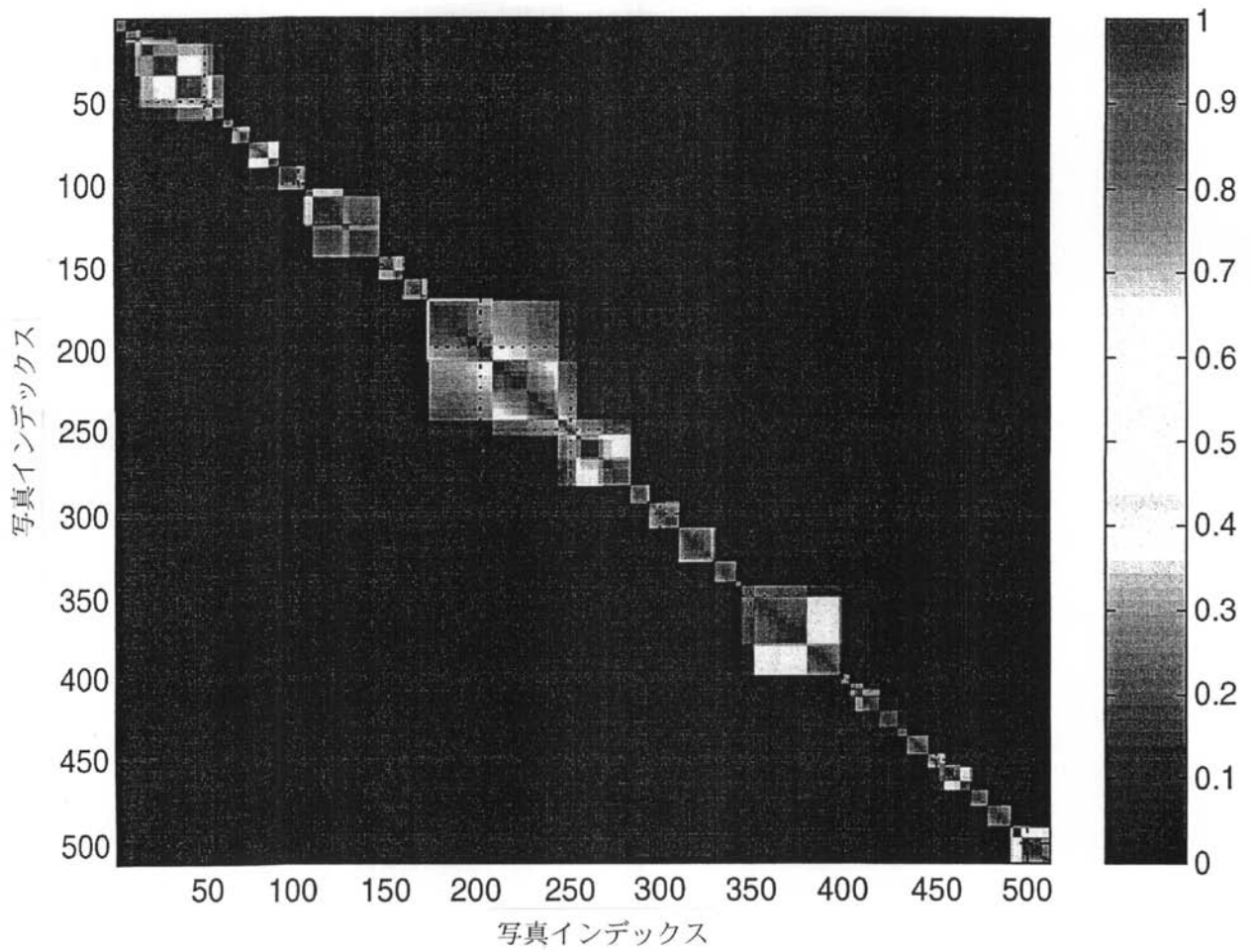
【 図 6 】



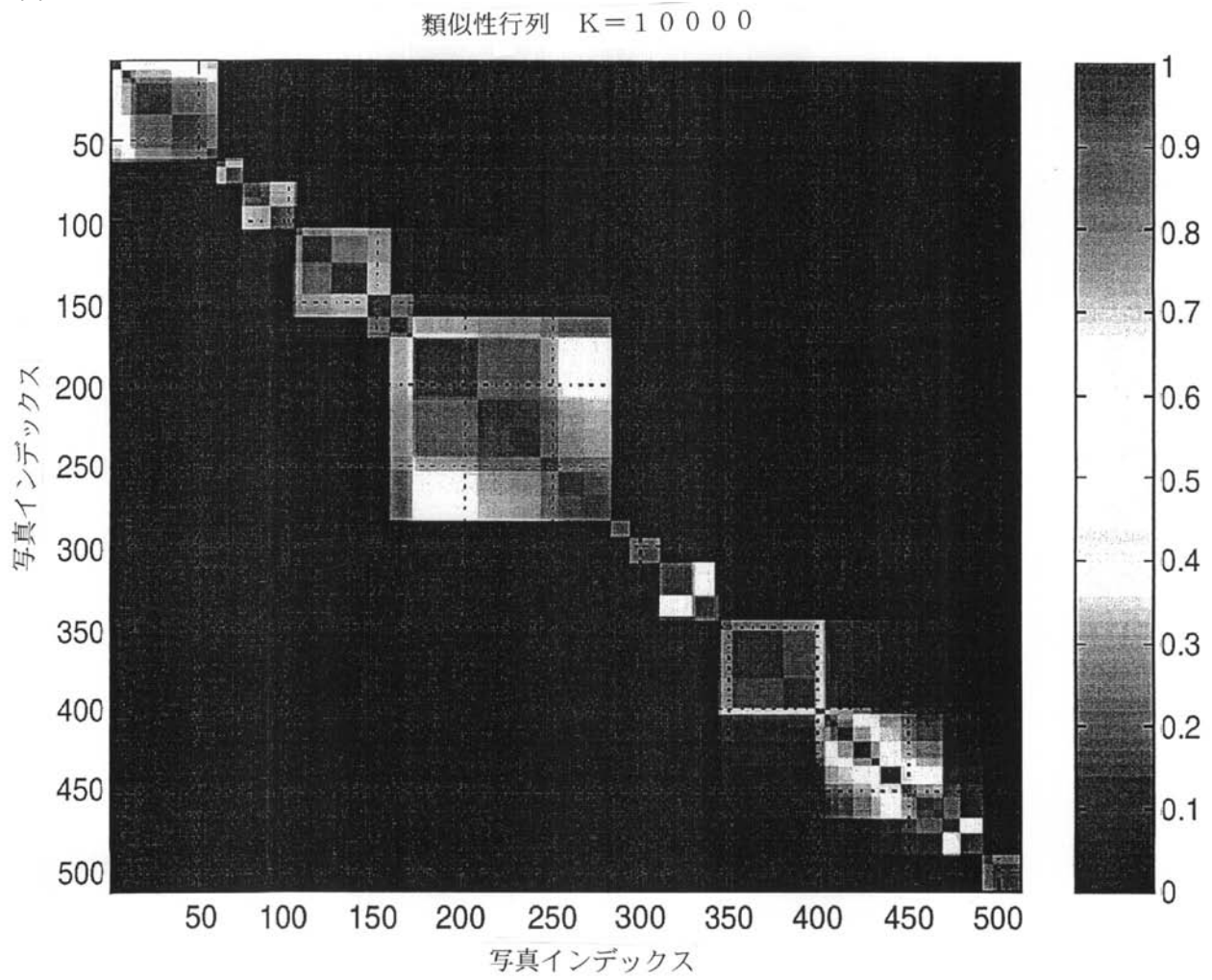
【 圖 1 8 】



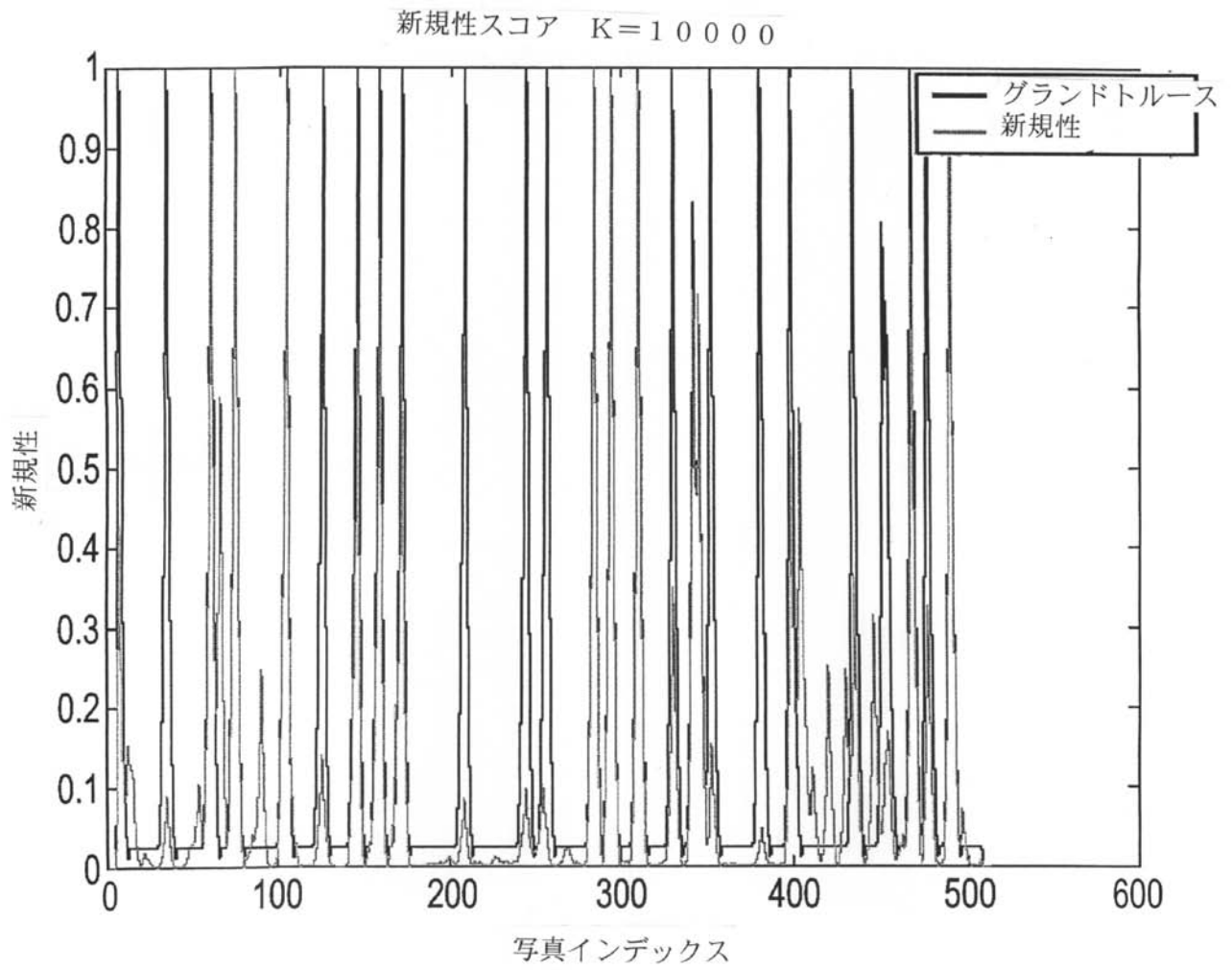
【図 5】

類似性行列 $K = 1000$ 

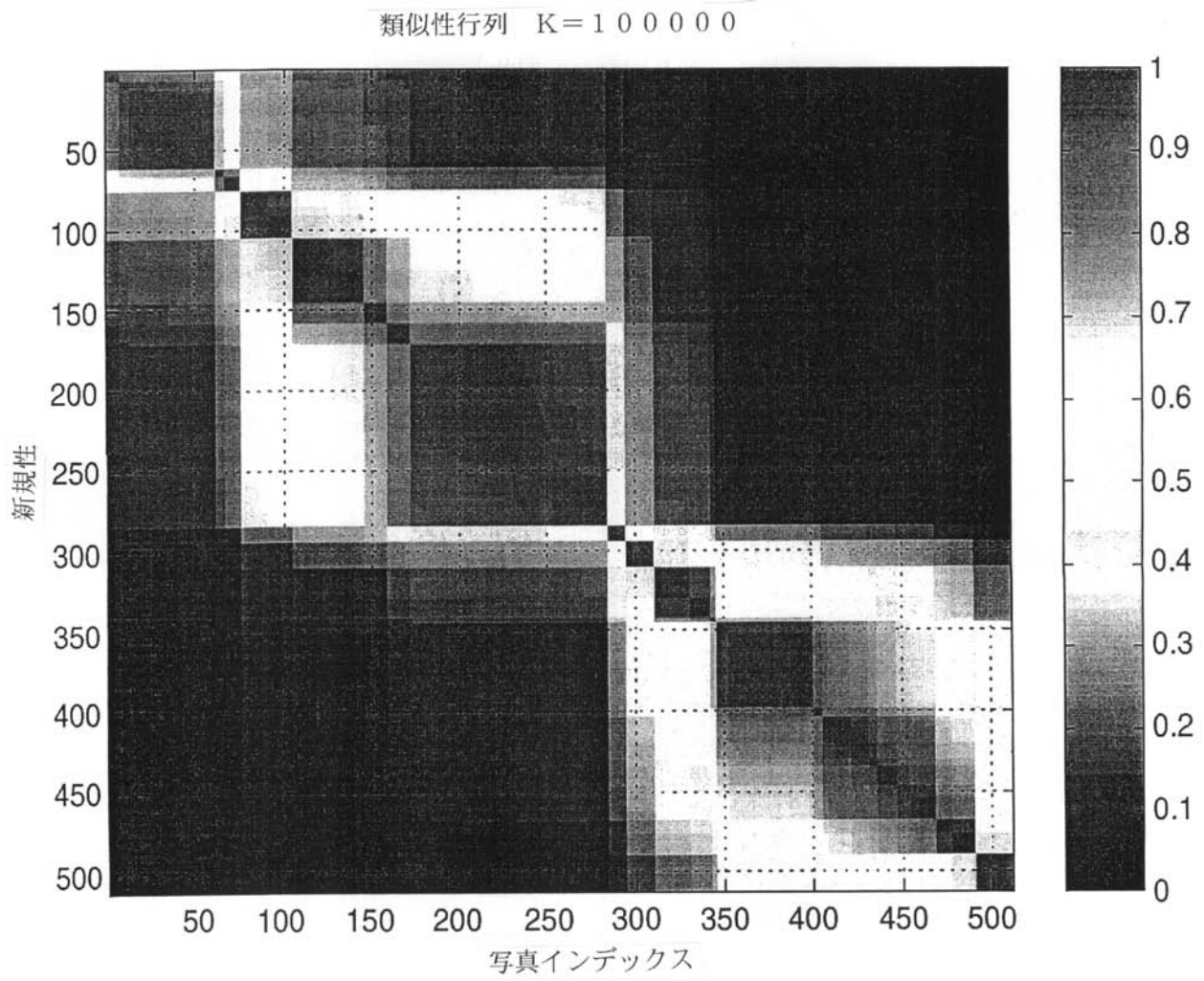
【図 7】



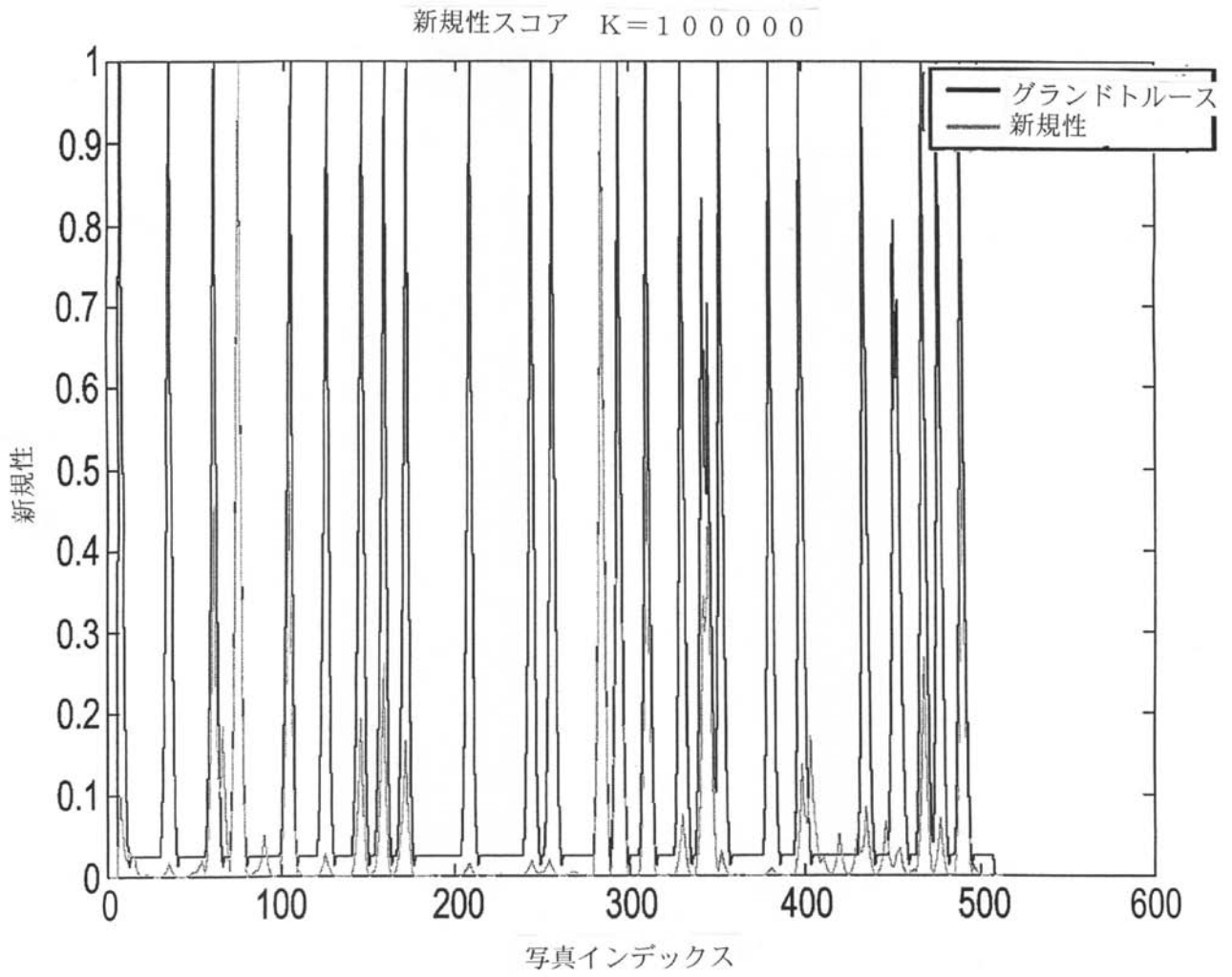
【図 8】



【図 9】

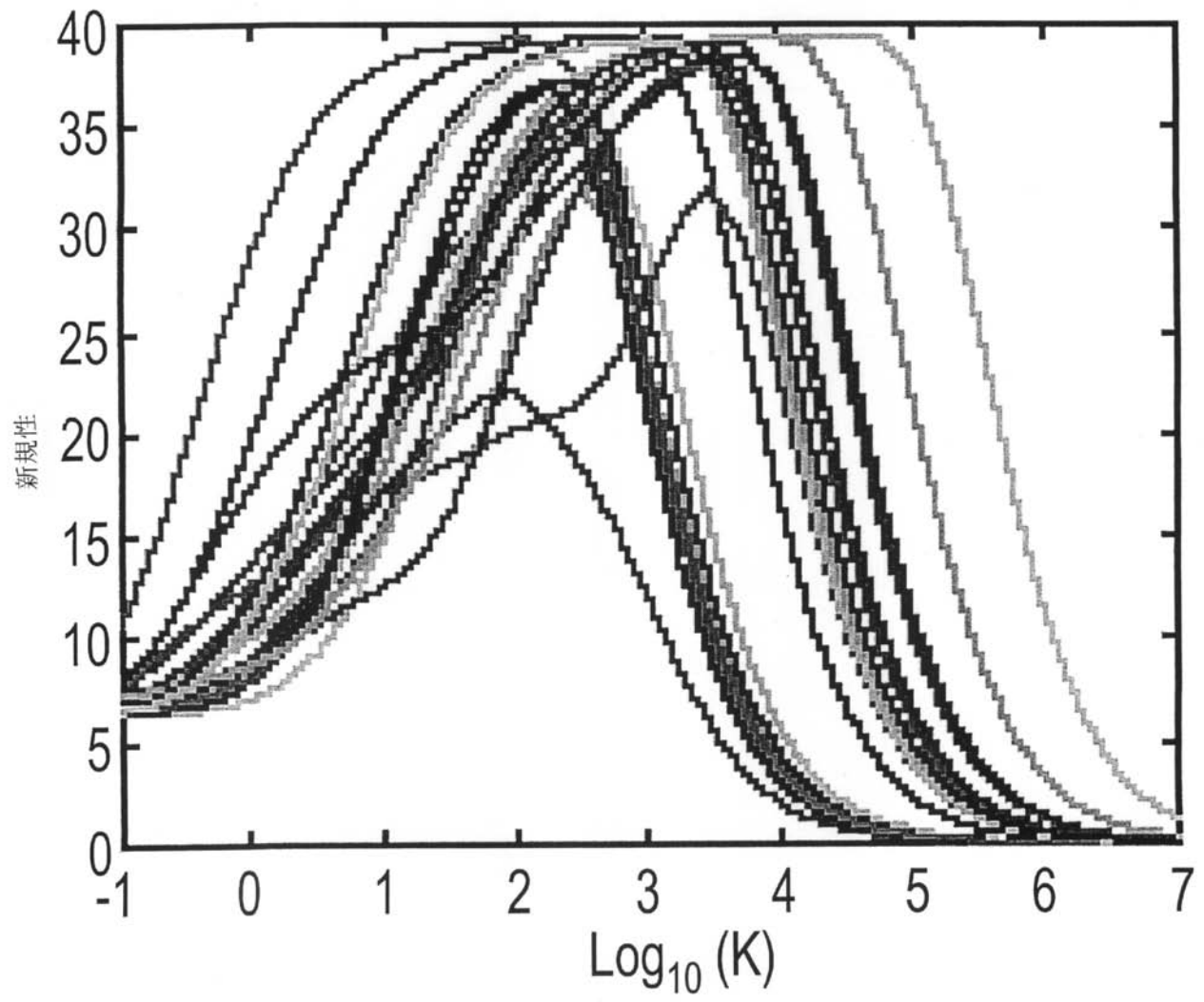


【図 10】

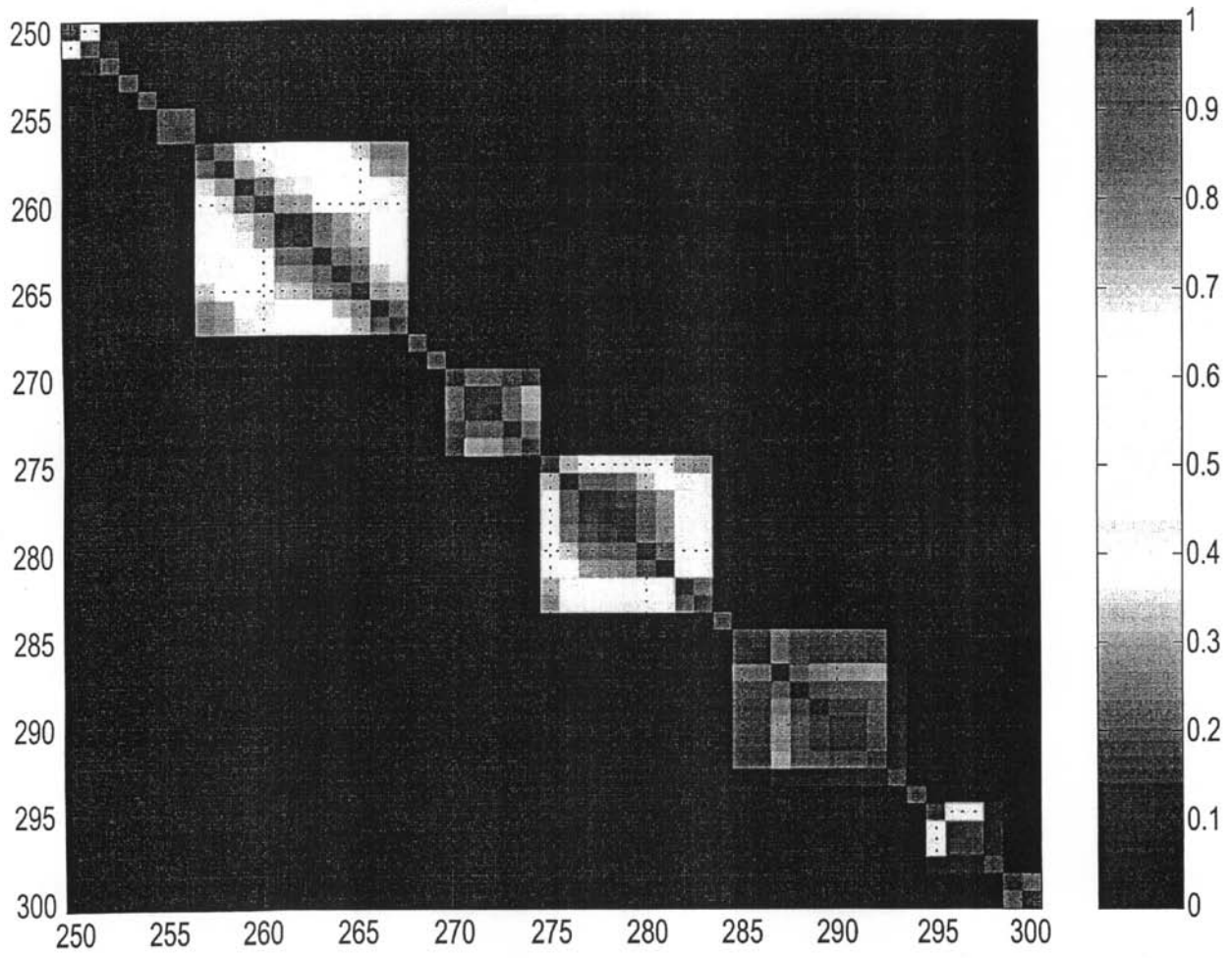


【図 11】

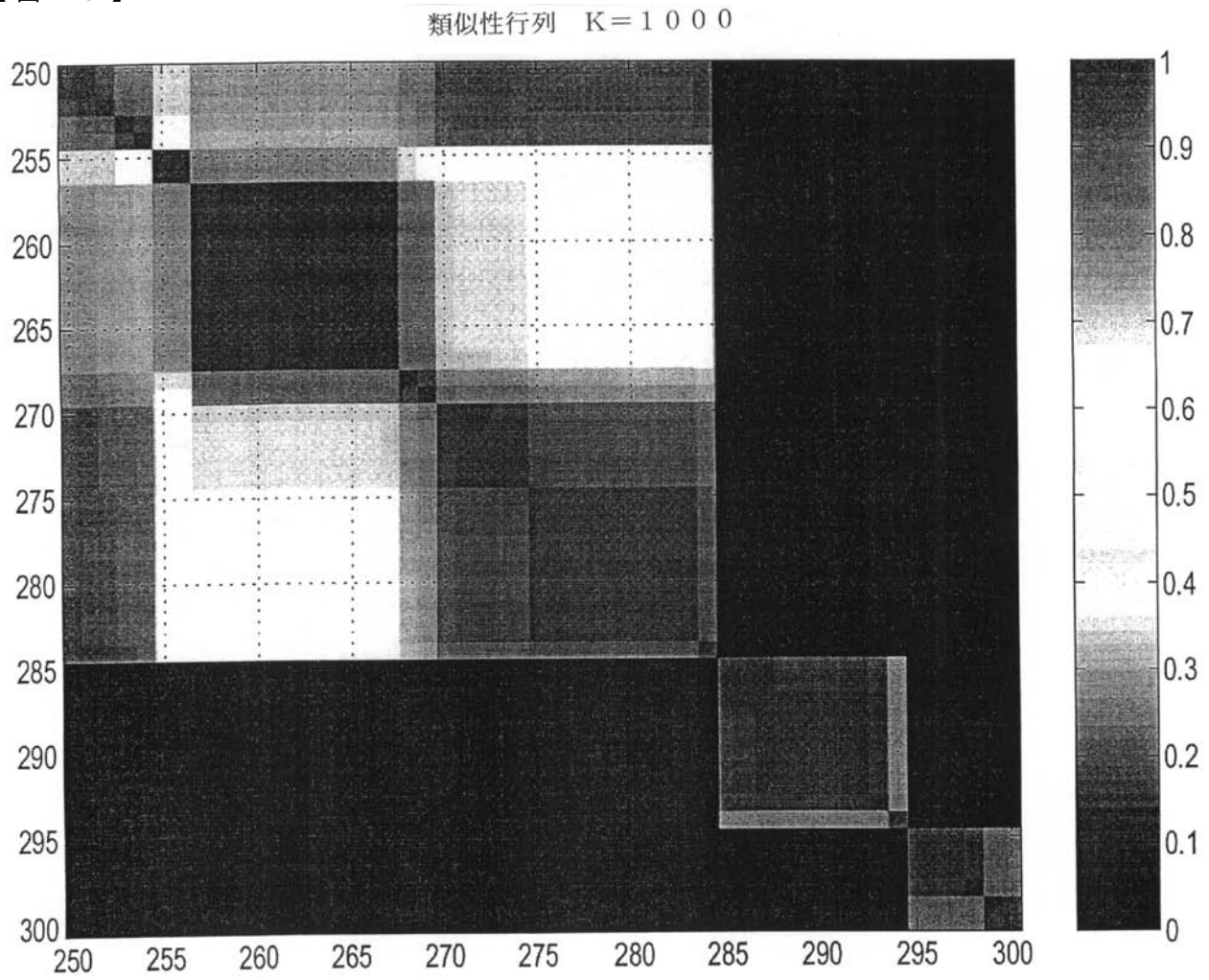
グランドトルースクラスタ境界についての新規性スコア



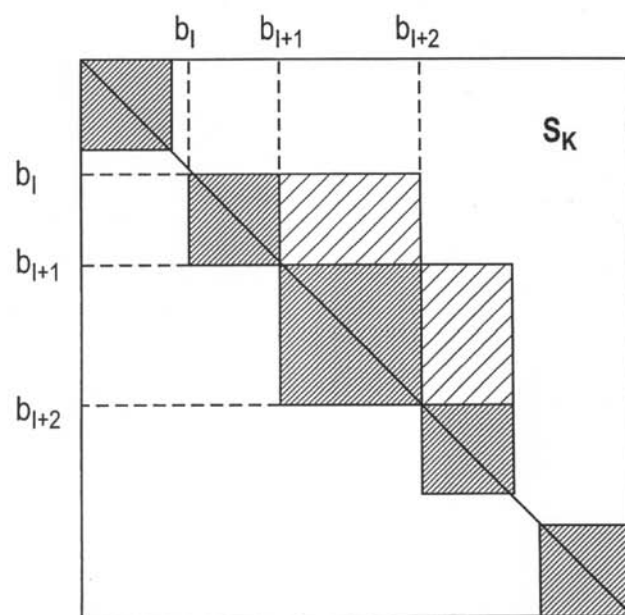
【図 1 2】

類似性行列 $K = 10$ 

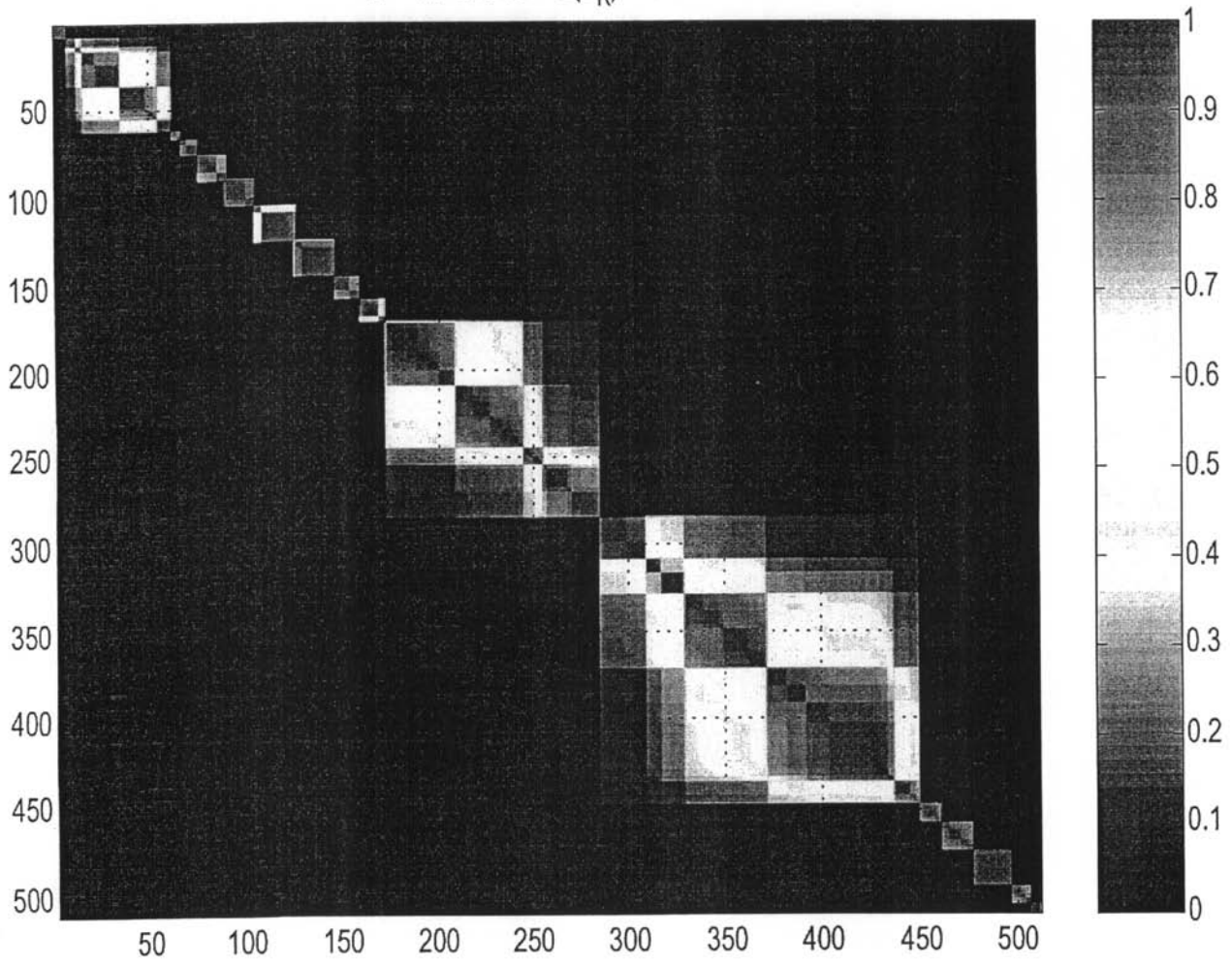
【図 1 3】



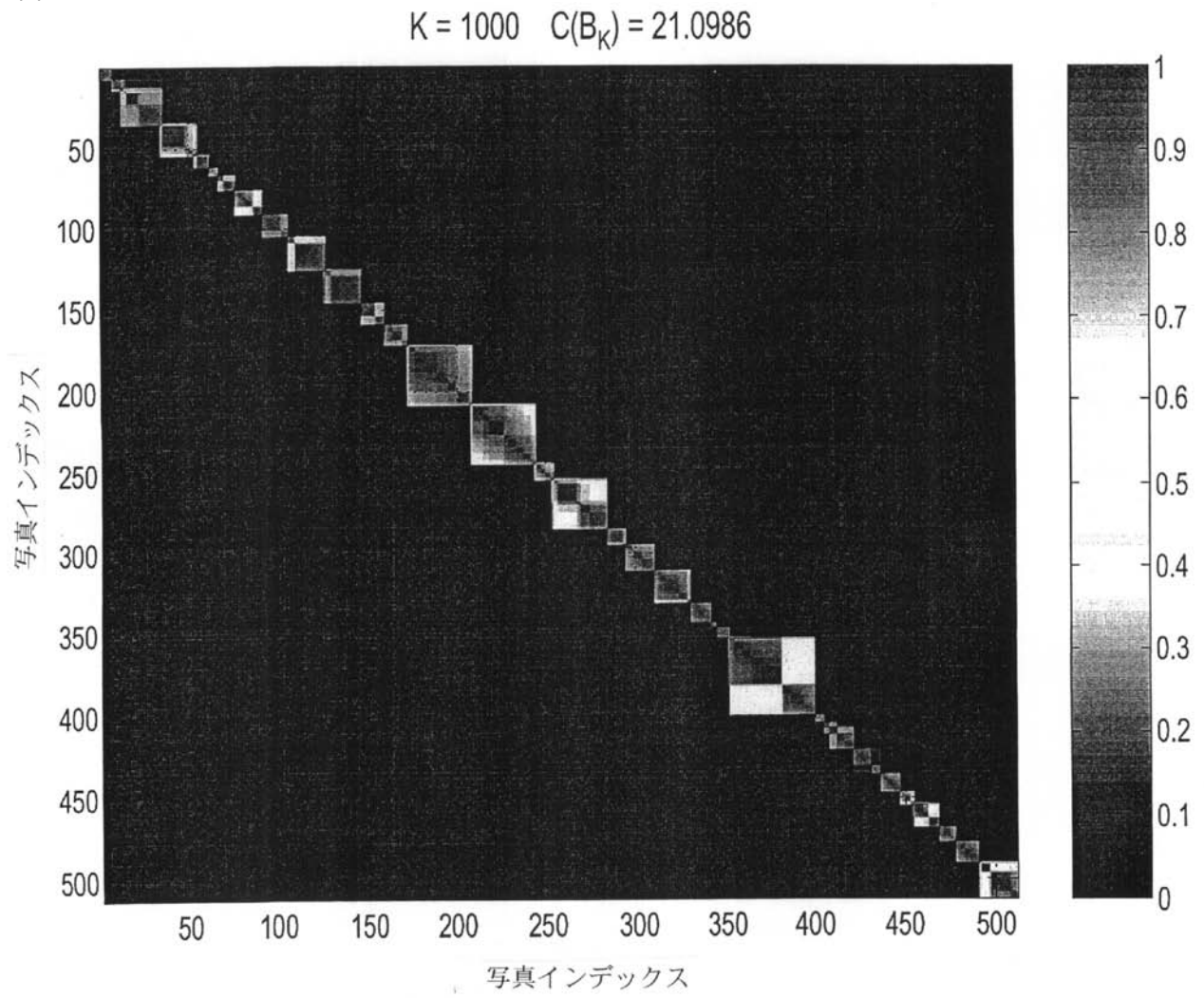
【図 1 4】



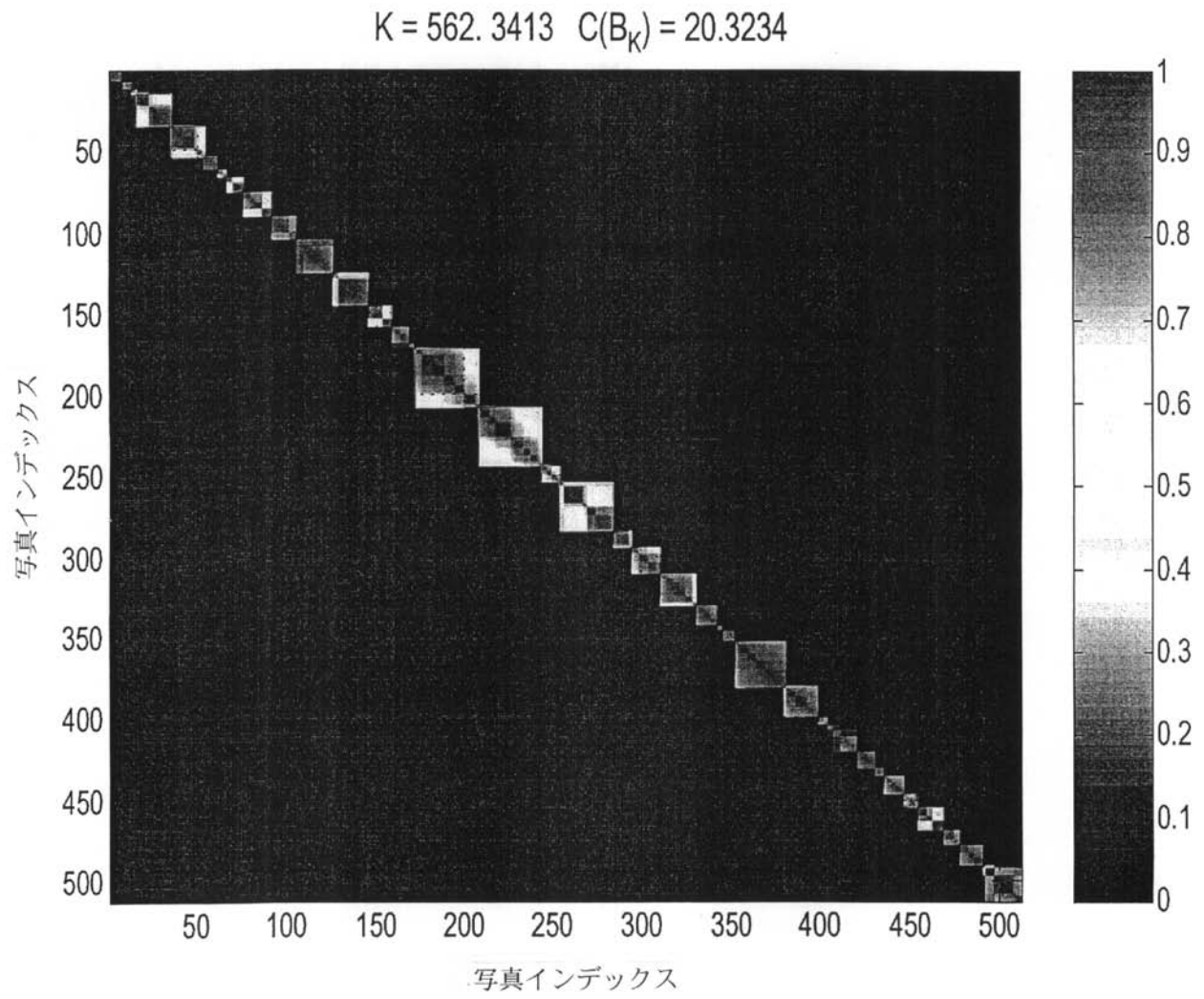
【 図 1 5 】

 $K = 1778.2794$ $C(B_K) = 11.7814$ 

【図 16】



【図 17】



フロントページの続き

(72)発明者 アンドレアス ガーゲンソン

アメリカ合衆国 9 4 0 2 5 カリフォルニア州 メンロ パーク ウェイバリー ストリート
2 1 0 ナンバー 4

F ターム(参考) 5B075 KK35 ND08 NR12 NS10 UU40