

(19) 日本国特許庁(JP)

(12) 特許公報(B2)

(11) 特許番号

特許第4874413号  
(P4874413)

(45) 発行日 平成24年2月15日(2012.2.15)

(24) 登録日 平成23年12月2日(2011.12.2)

(51) Int.Cl. F I  
**G06F 17/30 (2006.01)** G06F 17/30 350C  
 G06F 17/30 170G

請求項の数 3 (全 46 頁)

(21) 出願番号	特願2010-140222 (P2010-140222)	(73) 特許権者	596170170
(22) 出願日	平成22年6月21日 (2010.6.21)		ゼロックス コーポレイション
(62) 分割の表示	特願2000-16705 (P2000-16705) の分割		XEROX CORPORATION
原出願日	平成12年1月26日 (2000.1.26)		アメリカ合衆国、コネチカット州 068
(65) 公開番号	特開2010-250849 (P2010-250849A)		56、ノーウォーク、ピーオーボックス
(43) 公開日	平成22年11月4日 (2010.11.4)		4505、グローバー・アヴェニュー 4
審査請求日	平成22年6月21日 (2010.6.21)	(74) 代理人	100079049
(31) 優先権主張番号	117462		弁理士 中島 淳
(32) 優先日	平成11年1月26日 (1999.1.26)	(74) 代理人	100084995
(33) 優先権主張国	米国 (US)		弁理士 加藤 和詳
(31) 優先権主張番号	421416	(72) 発明者	フランシーン アール. チェン
(32) 優先日	平成11年10月19日 (1999.10.19)		アメリカ合衆国 94025 カリフォル
(33) 優先権主張国	米国 (US)		ニア州 メンロ パーク シャーマン ア
			ベニュー 975

最終頁に続く

(54) 【発明の名称】 オブジェクト間類似度計算方法

(57) 【特許請求の範囲】

【請求項1】

コンピュータが実行する、オブジェクトのコレクション内の2つのオブジェクトの間の類似度を計算するオブジェクト間類似度計算方法であって、

各オブジェクトは、少なくとも第1の特徴ベクトルと第2の特徴ベクトルに関連付けられ、

前記第1の特徴ベクトルと前記第2の特徴ベクトルはそれぞれ、複数の次元のベクトルであり、

前記第1の特徴ベクトルは、前記オブジェクトの第1の特徴を表し、

前記第2の特徴ベクトルは、前記オブジェクトの第2の特徴を表し、

前記第1の特徴は、テキスト特徴、URL特徴、インリンク特徴、及びアウトリンク特徴を含むマルチモードの特徴の第1のセット内の1つであり、

前記第2の特徴は、画像特徴であり、

前記オブジェクト間類似度計算方法は、

第1のオブジェクトの前記第1の特徴ベクトルと第2のオブジェクトの前記第1の特徴ベクトルとを識別するステップと、

前記第1のオブジェクトの前記第1の特徴ベクトルと前記第2のオブジェクトの前記第1の特徴ベクトルとの間の第1の距離メトリックを計算するステップと、

テキスト情報の参照なく、前記第1のオブジェクトの前記第2の特徴ベクトルと前記第2のオブジェクトの前記第2の特徴ベクトルとを識別するステップと、

前記第1のオブジェクトの前記第2の特徴ベクトルと前記第2のオブジェクトの前記第2の特徴ベクトルとの間の第2の距離メトリックを計算するステップと、

前記第1の距離メトリックと前記第2の距離メトリックとの合計を計算するステップと

を備えたオブジェクト間類似度計算方法。

【請求項2】

各オブジェクトは、ドキュメントのコレクション内の1つのドキュメントに対応する請求項1記載のオブジェクト間類似度計算方法。

【請求項3】

前記第1の特徴は、(i)前記テキスト特徴、(ii)URL特徴、(iii)インリンク特徴、及び(iv)アウトリンク特徴の少なくとも1つであり、

前記第1の距離メトリックは、前記第1のオブジェクトの前記第1の特徴ベクトルと前記第2のオブジェクトの前記第1の特徴ベクトルとの間のコサイン類似度である

請求項2記載のオブジェクト間類似度計算方法。

【発明の詳細な説明】

【技術分野】

【0001】

本発明は、オブジェクト間類似度計算方法、ドキュメント間類似度計算方法、及びユーザ特性間類似度計算方法にかかり、より詳細には、コンピュータが実行する、オブジェクトのコレクション内の2つのオブジェクトの間の類似度を計算するオブジェクト間類似度計算方法、コンピュータが実行する、ドキュメントのコレクション内の2つのドキュメントの間の類似度を計算するドキュメント間類似度計算方法、及び、コンピュータが実行する、ドキュメントコレクションのユーザ集団内の二人のユーザの特性間の類似度を計算するユーザ特性間類似度計算方法に関する。

【背景技術】

【0002】

コンピュータユーザは、求めているドキュメントコレクションを発見することが益々困難になってきている。理由は、そのようなコレクションのサイズが増加しているためである。例えば、インターネット上でワールドワイドウェブ(WWW)は、数百万の個々のページを含む。また、多くの会社の内部イントラネットは、しばしば数千ものドキュメントを含むレポジトリを含む。

【0003】

ワールドワイドウェブ上及びイントラネットレポジトリ中のドキュメントは、あまりうまくインデックス付けされていないということもしばしば真実である。その結果、特定のドキュメントのアイデンティティや位置や特徴が公知でない場合、そのような大きなコレクション(収集物)中で必要な情報を検出することは、望みの無い捜しものをするようなものである。

【0004】

ワールドワイドウェブは、インターネットを介して分散されたサーバ上に位置するドキュメント(大部分がテキストと画像である)の緩く相互リンクされたコレクションである。一般的に言えば、各ドキュメントは、例示のフォーム"`http://www.server.net/directory/file.html`"において、アドレス即ちユニフォームリソースロケータ(URL)を有する。その表記法において、"`http:`"は、そのドキュメントが引き渡されるべきプロトコルを指定し、この場合、"`HyperText Transport Protocol`(ハイパーテキスト転送プロトコル)である。"`www.server.net`"は、そのドキュメントが常駐するコンピュータやサーバの名前を指定し、"`directory`"は、そのドキュメントが常駐するサーバ上のディレクトリやフォルダを指し、"`file.html`"は、そのファイルのネームを指定する。

【0005】

ウェブ上の多くのドキュメントは、HTML（ハイパーテキストマークアップ言語）フォーマットであり、とりわけ、フォーマットिंगのそのドキュメントへの適用、外部コンテンツ（画像及び他のマルチメディアデータタイプのような）のそのドキュメント内への導入、及び他のドキュメントへの"ホットリンク"や"リンク"のそのドキュメント内への配置を可能とする。"ホットリンク"によって、ユーザは、ページ内の対象となるアイテムを選択することによって、簡単にウェブ上のドキュメント同士間をナビゲートすることが出来る。例えば、再グラフィック技術についてのウェブページは、ゼロックスコーポレートサイトへのホットリンクを持ち得る。ホットリンクを選択することによって（しばしば、マークされたワードや画像や領域をマウスのようなポインティングデバイスでクリックすることによって）、ユーザのウェブブラウザは、（通常、そのホットリンクに関連する、しばしばユーザには見えないURLを介して）ホットリンクに従って、異なるドキュメントを読み出すよう指示される。

10

**【0006】**

インターネット上の各及び全てのドキュメントに対するURLや、好ましいドキュメントの同様のコレクション中のこれらのドキュメントでさえユーザが思い出すことが期待できないことは明白である。従って、ナビゲーションの助けは、役に立つのみならず必要である。

**【0007】**

従って、ユーザがユーザのブックマークコレクションにはまだ表れていないインターネット（又は他の大きなネットワーク）上で情報を見つけることを希望する場合、ユーザは、しばしばその情報を検出するために"サーチエンジン"をオンにする。サーチエンジンは、インターネット上に格納されたコンテンツへのインデックスとして働く。

20

**【0008】**

サーチエンジンには、二つのカテゴリがあり、一つは、サブジェクトマターカテゴリの階層を形成するために解析され且つ使用されるドキュメント及びウェブサイトを含むカテゴリ（例えば、ヤフー（登録商標））であり、他は、用語のサーチ可能データベースを確立するためにウェブやドキュメントコレクションを"クロール"してページコンテンツ上のキーワードサーチを可能とするカテゴリ（とりわけ、AltaVista、Excite、及びInfoseek）である。

**【0009】**

また、ユーザにより或いは一つの好ましいドキュメントへの比較により提供される基準に基づきウェブサイト推奨を行うことが出来る推奨システムが知られている。

30

**【0010】**

前述のシステム及びサービスは、伝統的なサーチエンジンと推奨システム能力をある程度組合せるが、それらの何れもが全体としてインターネットよりも小さなグループのプリファランスを考慮すると共にサーチエンジンのような能力を提供するために現在適応できない。特に、コミュニティやクラスタベースの推奨（レコメンデーション）をインターネットや他のドキュメントのコレクションから未知のドキュメントを検索出来るシステムへ組合せることは有利である。

**【0011】**

従って、ドキュメントの大きなコレクション、即ちコーパスを処理する場合、それらのコンテンツ（内容）に基づき、これらのドキュメントをサーチし、ブラウジングし、検索し、そして見る事が出来る事が有用である。しかしながら、そのドキュメントの制限のために、これは多くの場合困難である。例えば、ワールドワイドウェブ上には、ドキュメントとファイルの一般的コレクションで利用できる多くの種類の情報がある。テキストファイル、テキストと画像の両方を含むHTML（ハイパーテキストマークアップ言語）、画像、サウンドファイル、マルチメディアファイル、及び他のタイプのコンテンツがある。

40

**【0012】**

画像を容易にブラウジングし検索するために、コレクション中の各画像は、理想的には

50

、その画像中のオブジェクト及びその画像の記述を含む記述情報でラベル付けされるべきである。しかしながら、ウェブ上のオブジェクトのような、画像の無制限のコレクション中のオブジェクトの識別は困難なタスクである。オブジェクトの自動的識別方法は、一般的に、マシーンパーツのような特定のドメインに制限される。人間に各画像を識別させることは、ウェブ上におけるように、面倒な仕事であり、しばしば不可能である。

【0013】

情報検索における多くのサーチは、コンテンツに基づくテキストドキュメントの検索や、可視特徴に基づく画像ドキュメントに基づく検索に焦点が当てられている。更に、ウェブ及び会社のイントラネット上の情報の急激な増加で、ユーザは、特定の情報をサーチする時、ヒント攻めに遭う。実際に求められる物を発見するために結果を分類するタスクは、しばしば冗長であり、時間がかかる。最近、多くのサーチエンジンは、ユーザがメタデータ（例えば、Hotbot、Infoseek）を介して伝統的なキーワードエントリから質問を拡張することを可能とする機能性を追加している。メタデータは、言語、日付、サイトの位置、或いは画像、ビデオ又はオーディオ等のモダリティが存在するか否かのような、種々のフォームを取ることが出来る。

10

【0014】

しかしながら、最近、検索のために使用マルチモード特徴についてのあるリサーチがある。コレクション中のドキュメントのマルチモード特徴並びにユーザのブラウジングの習癖の間の類似度に基づいて、ユーザが希望の情報を検出できる幾つかのアプローチがここで示される。

20

【0015】

ユーザが画像とその画像と関連するテキストの両方、及び使用のような、ドキュメントに関連する他のタイプの情報を使用してサーチを反復的に狭めるドキュメントブラウジング及び検索へのアプローチがここで記述される。テキスト、画像特徴及び使用のような、異なるタイプの情報は、“モダリティ”と呼ばれる。オブジェクトに関連する幾つかのモダリティからのデータを有するオブジェクトのグループ化は、ここでは、マルチモードクラスタリングである。

【0016】

画像を囲む或いはそれと関連するテキストは、しばしばそのコンテキストの指示を提供する。ここで提案される方法は、（ここで述べられる例示のケースでは、画像の）ブラウジング及び検索を実行するために、テキストと画像の特徴のようなマルチモード情報の使用を可能とする。この方法は、より一般的には、コレクションの要素（例えば、ドキュメント、フレーズ或いは画像）が複数の特性、即ち特徴によって述べられ得る他の用途へ適用出来る。

30

【0017】

サーチ及びブラウジングにおける複数の特徴の使用での一つの困難なことは、異なる特徴から情報を組合せることである。これは、一般には、ユーザによって設定され得る、各特徴（通常、カラーヒストグラム、テクスチャー及び形状）と関連する重みを有することによって画像検索タスクで扱われる。重みを変更される毎に、新たなサーチが実行されなければならない。しかしながら、異質セットのマルチモード特徴の使用において、しばしば異なる特徴の重要度へ重みを割り当てるのが困難である。メタデータを使用するシステムにおいて、有限で離散的な値を有し、特定の値を含むか排除するブールシステム（系）が使用され得る。この概念を離散的リードではないマルチモード特徴へ拡張することは、特徴の組合せ方の問題を一層悪化する。

40

【先行技術文献】

【特許文献】

【0018】

【特許文献1】特開平10-27125号公報

【発明の開示】

【発明が解決しようとする課題】

50

## 【0019】

従って、種々のコンテキスト及び用途において、フレキシブルにマルチモード情報を処理できるシステムが必要である。照会を実行すると共に、直接テキストコンテンツ以外の特徴、即ち、画像特徴及び間接特徴によってサーチ結果を詳細化し調節できることは有用である。また、個人ユーザがアクセスするドキュメントの特徴を介してその個人の情報アクセス習性を追跡でき、それによってユーザが類似クラスタへ割り当てられる推奨システムを可能とすることは有用である。

## 【課題を解決するための手段】

## 【0020】

この開示は、マルチモードブラウジングとクラスタリングに対するフレームワークを記述し、ドキュメントのコレクション中のコンテンツのブラウジング、サーチ、検索及び推奨を向上するためにそのフレームワークを有利に使用するシステムを記述する。

10

## 【0021】

大きなデータセットのクラスタリングは、探索データ解析、可視化、統計的汎用化、及び推奨システムにとって重要である。大部分のクラスタリングアルゴリズムは、オブジェクト間の類似度に依存する。このプロトコルは、データ表示モードとマルチモードデータに対する関連する類似度メトリックを記述する。このアプローチは、各オブジェクトがそれと関連する幾つかの異なるタイプの情報を有するデータセットと関連し、モダリティと呼ばれる。そのようなデータセットの例は、ワールドワイドウェブ（ここで、モダリティは、テキスト、インリンク、アウトリンク、画像特徴、テキストジャンル、等）を含む。

20

## 【0022】

本発明の主な特徴は、その新規のデータ表示モードにある。各ドキュメント内の各モダリティは、ここでは、 $n$ -次元ベクトルによって記述され、コレクション内の関連の量的解析を容易にする。

## 【0023】

本発明の一用途は、情報をブラウジング及び検索するために異なる空間（即ち、異なるモダリティ）におけるドキュメント特徴を連続して使用するための方法が記述される。この方法の一実施の形態は、画像のブラウジング及び検索のために画像とテキスト特徴を使用するが、本方法は、一般に、あらゆるセットの個々の特徴に適用される。本方法は、ユーザが対象となる項目（アイテム）を指定できる複数の方法の利点を有する。例えば、画像において、そのテキストと画像のモダリティから特徴がその画像を記述するために使用され得る。この方法は、米国特許第5,442,778号、及び1992年の第15周年 Int'l SIGIR'92の会報のD. Cutting、ドクターKarger, J.O.、Pedersen 及び J.W. Tukey による "Scatter/Gather: A Cluster-based approach to browsing largedocument collections" ("Scatter/Gather") に開示された方法と類似する。この "Scatter/Gather" では、クラスタの選択及びその後続く選択されたクラスタの再クラスタリングが反復的に行われる。それは、少なくとも二つのこと、即ち、各クラスタリングが異なる特徴（例えば、周りのテキスト、画像URL、画像カラーヒストグラム、その周りのジャンル）に実行できること、及び "マップ" 機能が指定された特徴に関する最も類似するクラスタを識別すること、に関して、Scatter/Gatherパラダイムを拡張する。後者の機能は、特徴値の欠落に起因して除外されている追加の類似画像の識別を可能とする。画像クラスタは、各クラスタから少数の代表的画像を選択することによって表示される。

30

40

## 【0024】

本発明の他の用途では、異なるモダリティにおける種々のドキュメント特徴は、全体の類似度を表すクラスタを形成するために、適切に重み付けされ組み合わせられる。

## 【0025】

また、本発明の種々の代替の実施の形態は、一つ又はそれより多くの特徴に従ってユーザ及びドキュメントのクラスタリング、ユーザクラスタの従来のブラウジング動作に基づくドキュメントの推奨、及び図形及びテキストでのドキュメント又はユーザのクラスタの可視表示を可能とする。

50

## 【 0 0 2 6 】

最初に、ベクトル空間におけるユーザ及びドキュメントの表示を行い、ウェブ画像のコレクション及びHTMLページ上の関連テキストのブラウジング及び検索を実行するシステムが記述される。ブラウジングは、ユーザによる情報のコレクション又はコーパスの対象となる部分の検出を助けるために、このコーパスに良好にマッチする照合を公式化する必要無く、検索と組み合わせられる。画像を囲むテキスト及び幾つかの単純な画像特徴の形態で、マルチモード情報がこの処理で使用される。このシステムを使用して、ユーザは、コレクションを少数の対象となる要素に徐々に狭める。これは、マルチモード特徴を使用するために拡張されることを除いてテキストブラウジングのために開発されたScatter/Gatherシステムと類似する。上述のように、幾つかのコレクションの要素は、幾つかの特徴に対して未知の又は未定義の値を有してもよい。一つの方法は、これらの要素を結果としてのセットへ結合する方法が示される。また、この方法は、サーチが二つのクラスタ同士の境界近くの空間の一部へ狭められるケースを扱う方法を提供する。多くの例が提供される。

10

## 【 0 0 2 7 】

種々のメタデータフィールドと有するデータベースと類似して、本発明のコレクションにおけるドキュメントは、多くの異なる特徴、即ち、多くが非構成ドキュメントのコンテンツから導出される(多分、非直交)"次元"によって特徴付けられる。

## 【 0 0 2 8 】

マルチモード特徴は、ユーザ情報、テキストジャンル、又は画像解析等の多くのフォームを取り得る。本発明において使用される特徴は、メタデータが一般的に手作業で割り当てられる現在の画像サーチシステムとは異なり、データ(例えば、テキストと画像)及びそのコンテキストから導出される且つ自動的又は半自動的に割り当てられるメタデータのフォームであると考えられ得る。表1は、幾つかの可能な特徴(それらの全ては、より詳細に後述される)を示す。種々の他の特徴及びモダリティもまた本発明において使用でき、且つ表1の特徴は、例示に過ぎない事が理解される。

20

## 【表1】

特徴	モダリティ
Text Vector	テキスト
Subject	テキスト
URLs	テキスト
Inlinks	ハイパーリンク
Outlinks	ハイパーリンク
Genre	ジャンル
Page Usage	ユーザ情報
Color Histogram	画像
Complexity	画像

30

40

## 【 0 0 2 9 】

ここでは、豊富な"マルチモード"特徴を組合せてユーザに情報のニーズを満足させる方法が提供される。一方では、これは、アドホック検索(画像に適用される)を含み、ユーザのニーズに関連する情報への簡単で迅速なアクセスを提供する。他方では、これは、解析されるドキュメントコレクションとそれらのユーザを含む。共通のシナリオは、ワールドワイドウェブであり、それは、多くの大量のドキュメントコレクションに典型的な種類の非構成ドキュメントから成る。

## 【 0 0 3 0 】

従って、この明細書は、ウェブ画像のコレクション及びHTMLページ上の関連するテ

50

キストに対する情報アクセスの方法を提供する。この方法は、画像及びそれらの画像に関連するドキュメント又はドキュメント領域のブラウジング及び検索を実行するために、テキストと画像の特徴のようなマルチモード情報の使用を可能とする。記述されたアプローチにおいて、画像の内容を指示する、その画像を囲む又はそれと関連するテキストから導出されるテキスト特徴は、画像特徴と共に使用される。このアプローチの新規性は、ユーザに見えるテキスト及び画像特徴を作る方法に依存し、ユーザが対象となるユーザのサーチを連続的に狭めることを可能とする。これは、特にユーザがコーパスに良好にマッチする照合を公式化することが困難な時、特にコーパスで使用されるボキャブラリや画像記述子が未知である場合にウェブのような一般的でない又は異質コーパスと共に働く時に有用である。

10

**【0031】**

ここで開示される方法は、ドキュメント（とユーザ）特徴が多次元ベクトル空間内に埋め込まれている有利なデータ表示モデルを前提とする。このデータ表示モデルは、一貫性があり対称的類似度尺度の使用を容易とし、それは、以下で詳細に説明される。ここで記述されるデータ表示と類似度モデルによって、ユーザ（即ち、コレクション使用データ）によってアクセスされるドキュメントのコンテンツと特徴に基づいて、ユーザとユーザのクラスタを表示することが出来、それによって、それらの類似度に従ってユーザをクラスタリングする能力を改良する。

**【0032】**

更に、マルチモードユーザクラスタに基づく推奨システムは、後述されるように、マルチモードコレクション使用データのコレクションで可能とされる。一セットのクラスタは、ユーザのトレーニングセットから誘導される。推奨を望むユーザは、最も近いクラスタへ割り当てられ、このクラスタの好ましいドキュメントがユーザへ推奨される。

20

**【0033】**

最後に、ここでの開示は、ドキュメントのクラスタ及びユーザのクラスタを可視表示する改良された方法を示す。ドキュメントがしばしば階層的に格納されて、階層可視表示を可能とするが、それは、ユーザに対して通常は真であるとは限らない。従って、本発明は、適切なユーザによってアクセス又はアクセスされ得るドキュメントの階層的ビューを介してユーザデータを見るのが可能となる。ドキュメント及びドキュメントのクラスタは、クラスタの"顕著な次元（ディメンション）"を介して類似的に且つテキスト的に可視化され得る。

30

**【0034】**

画像検索においてクラスタリングの使用は新しいことではないが、クラスタリングは、データベース母集団化ステージの間に人を助けるために、或いは画像をオフラインでクラスタリングして照合同士間の距離サーチがクラスタ内で実行されるように、事前処理のために使用されている。本発明において、反復的なクラスタリング及びクラスタサブセットの選択は、対象となる画像の識別を助けることが出来る。クラスタリングは、反復サーチ及び提示のために使用され、適合フィードバックは、クラスタのユーザの選択において内在する。また、ユーザが個々の画像ではなくてクラスタを扱っているので、フィードバックステップは、実行がより容易である。

40

**【0035】**

ここで述べられるマルチモードクラスタリングの種々のフォームは、情報アクセスのため、即ち、ドキュメントを見つけるためにコレクションをブラウジングするため、ユーザにとって新しいコレクションを理解するため、及び"なにも検出されない"場合（クラスタリングは、ユーザがユーザの照合をコレクションに適するボキャブラリに公式化することによってユーザの照合を再公式化することを助ける事が出来る）を扱うために、使用され得る。

**【発明を実施するための最良の形態】****【0036】**

概説的に上記されたように、コレクション中のドキュメントを効率的にブラウジングし

50

サーチするために本発明のシステム及び方法の能力は、矛盾の無いデータ表示モデルの存在に大きく依存する。具体的には、ドキュメント同士の間の定量的類似度メトリックを定義するために、マルチ次元ベクトル空間にドキュメントをマッピングすることが有用であることが判った。従って、ここで記述されるアプローチは、全てのモダリティに対してデータ表示モデルを定義し、ここで各ドキュメントは、 $R^n$ として表される。このモデルは、図1に最も良く図示されている。

【0037】

図1に示されるように、コレクション120から選択された各ドキュメント（例えば、HTMLドキュメント110）は、各モダリティ（例えば、テキストベクトル114とURLベクトル116）毎に一つ、一セットの特徴ベクトル112へマッピングする。

10

【0038】

特徴ベクトル112は、ドキュメントコレクション120と通信ネットワーク124（インターネットや会社のイントラネットのような）の両方へのアクセスを有するプロセッサ122によって計算される。本発明の一実施の形態において、コレクション120は、ネットワーク124へも接続された一つ又はそれより多くのサーバによってホストされる。各ドキュメント毎の特徴ベクトル112は、データベース126へ格納され、そこでそれらの特徴ベクトルは、対応するドキュメントと照合される。ネットワークへ接続された複数のユーザ端末128、130及び132は、システムをアクセスするために使用される。

【0039】

20

これらの特徴ベクトルは、ドキュメントが初めにそのコレクション120に追加された時又はその後、システムによって発生される。本発明のこの実施の形態において、コレクション120は、本発明のシステムによって今までに処理された全ての既知のドキュメントから成ることが観察されるべきである。しかしながら、サーチエンジンの照合の結果に対して即時にコレクションを発生することもまた可能である。次に、極端に大きなグループのドキュメント（ワールドワイドウェブのような）に対してより実践的である、このアプローチは、オリジナルサーチ結果を編成、ブラウジング、ビュー或いは扱うために使用できる。

【0040】

ドキュメントをコレクション120へ追加するこの動作は、図2に示されるように、実行される。初めに、新たなドキュメントが検出される（ステップ210）。そのドキュメントは、特徴ベクトル112を計算するために処理される（ステップ212）。次にそのドキュメントは、本発明で利用できるコレクション即ちコーパスへ追加され得る（ステップ214）。もはやドキュメントが無い時（ステップ216）、処理は終了する（ステップ218）。そうでない場合、他のドキュメントが検出され（ステップ210）、処理が繰り返される。

30

【0041】

そのシステムのここで示される好ましい演算ベクトルは、八つの可能なドキュメントベクトル特徴、即ち、テキストコンテンツ、ドキュメントリンク、インリンク、アウトリンク、テキストジャンル、画像カラーヒストグラム、及び画像複雑さ、を使用できる。リストアップされた特徴の最初の二つは、テキストベースであり、インリンクとアウトリンクは、ハイパーリンクベースであり、テキストジャンルは、確率ベースであり、そして最後の二つの特徴（画像カラーヒストグラムと画像複雑さ）は、画像ベースである。これらの特徴は、それらの単純性と理解可能性のために、本発明と共に使用するために選択された。選択された特徴は、情報アクセスにおいて画像とテキストモダリティを使用及び組合せるための本発明方法を図示するように働く。しかしながら、多くの他のドキュメントメトリック（異なる画像領域に対するローカルカラーヒストグラム、画像セグメンテーション、及びテキスト特徴から2, 3ではあるがネームまで）もまた可能であり、本発明のシステム又は方法の範囲内で開発可能であることが理解される。

40

【0042】

50

本発明の一実施の形態において、これらの特徴ベクトルが図3に示されるように導出される。新たなドキュメント（それらは、テキストドキュメント、画像又は他のタイプの情報であってもよい）コンテンツが分離された（ステップ310）後、本発明の方法は、特徴ベクトルを導出するために、種々の情報源を使用する。テキストは、ドキュメントから抽出され（ステップ312）、対応するテキストベクトルを生成し（ステップ314）、対応するURLベクトルを生成する（ステップ316）ために使用される。

【0043】

他方（同時に或いはそれに続いて）、全てのアウトリンク（他の場所を指すドキュメント内のハイパーリンク）が抽出され（ステップ318）、対応するアウトリンクベクトルを生成する（ステップ320）ために使用される。インリンク（主題のドキュメントを指すコレクション内のドキュメント）が抽出され（ステップ322）、対応するインリンクベクトルを生成する（ステップ324）ために使用される。テキストジャンルが識別され（ステップ326）、対応するジャンルベクトルを生成する（ステップ328）ために使用される。

【0044】

その新たなドキュメントが少なくとも一つの画像である或いはそれを含む場合、次に、カラーが画像から抽出され（ステップ330）、対応するカラーヒストグラムベクトルを生成する（ステップ332）ために使用される。単一のカラー（又は一セットの類似のカラー）の水平及び垂直ランもまた画像から抽出され（ステップ334）、カラー複雑さベクトルを生成する（ステップ336）ために使用される。

【0045】

最後に、ドキュメントに対する参照が使用ログから抽出され（ステップ338）、ユーザのページアクセスベクトルを更新する（ステップ340）ために使用される。

【0046】

次に、コンテンツベクトルの全てがデータベースに格納される（ステップ342）。

【0047】

上述の異なる特徴ベクトルタイプを計算するための本発明の方法が以下に詳細に説明される。

【0048】

しかしながら、ある特徴を有するドキュメントを既存のコレクションに追加することは、コレクション中の全てのドキュメントに対する全セットの特徴ベクトルを訂正することが必要である場合があることを理解すべきである。例えば、特異なワードを含むドキュメントを追加することは、そのワードが余分の用語を各ドキュメントのテキストベクトルへ追加することを必要とするので、そのコレクション中の全てのドキュメントに対するテキストベクトルに影響を及ぼす。従って、実質的に大きなグループのドキュメント中のコレクションを更新することは、新たなドキュメントが利用出来るようになる毎に増分的に更新するよりも演算的により効率的である。このような考察及びそのセットのベクトルを演算的に最適化する方法が導入されるが、その詳細は、本発明にとって重要ではない。

【0049】

本発明の一実施の形態において、各特徴が別々に使用され、最も適切な距離メトリックが各特徴に適用され得る。本発明の他の実施の形態において、それらの特徴は、そのドキュメントを表す単一のコンテンツテキストベクトルへ組み合わせられ、そして単一の距離メトリックがそれらのドキュメントをクラスタリングし比較するために使用される。これらの他の実施の形態は、以下により詳細に記述される。

【0050】

ドキュメント情報のベクトル空間表示  
各タイプの特徴ベクトルの計算は、以下でより詳細に記述される。しかしながら、以下に見られるように、幾つかの一般的な特徴が全ての表示に当てはまる。

【0051】

テキスト特徴は、図4に示されるように計算される。テキスト特徴は、用語ベクトルで

10

20

30

40

50

あり、そこで、そのベクトルの要素は、ドキュメント自体で使用される用語を表す。本発明のこの実施の形態において、全てがテキスト又はHTMLドキュメント（又はテキストを実際を含む他のドキュメントタイプ）に対して、テキストベクトルは、ドキュメントの全体のテキストコンテンツに基づく。そのドキュメントが画像（又は実際のテキストを含まない他のタイプのドキュメント）の場合、テキストベクトルを公式化するために使用されるテキストは、"ホスト"HTMLページ中の画像を囲むテキストから導出される。この囲むテキストの範囲は、画像位置の前又は後の800文字に制限される。水平ルール、ヘディング又は他の画像がその制限に達する前に生じると、その範囲は、そのルール、ヘディング又は画像で終わる。"ストップリスト"は、冠詞、前置詞及び接続詞のような殆ど内容の無い共通の用語のインデックス付けを防止するために使用される。

10

## 【0052】

従って、ここで記述される本発明の目的のために、テキストドキュメント、画像ドキュメント及びマルチメディアドキュメントは、全て総称的用語"ドキュメント"の特別のケースであり、これらの特別なケースの各々に対して、ここで記述されるモダリティの幾つか又は全てが適用出来る。例えば、上述されたように、画像は、必ずテキストを含む必要は無いが、それらを指すハイパーテキスト又はURL中にテキストによって記述される。テキストを含む画像（ファクシミリビットマップ）は、既知のドキュメント画像デコーディング技術を介して抽出されたそれらのテキストを有してもよい。同様に、オーディオファイルは、またハイパーリンク及びURL中のテキストにより参照されることが出来、且つ既知のスピーチ認識アルゴリズムを介して抽出可能なテキストを含むことが出来る。幾つ

20

## 【0053】

上で提案されたように、ここで記述されるベクトル空間モデルにおいて、各テキストドキュメント $d$ （又は、あらゆる種類のテキストを含むドキュメント）が本発明によって $R^{n_t}$ に埋め込まれるが（ベクトル空間が $n_t$ 次元を含み、そこで各次元が実数で表される）、ここで $n_t$ は、コレクション中の特異なワードの全数（ $n_t$ はテキスト要素の数を表す）である。ベクトル空間への埋め込みは、以下のように定義される。

$$t(d)_i = tf_{di} icf_i$$

30

ここで、 $d$ は、特定のドキュメントであり、 $i$ は、ワードのインデックスであり、 $t(d)_i$ は、ベクトル $t(d)$ の要素 $i$ である。トークン頻度重み( $tf$ )と逆コンテキスト頻度重み( $icf$ )は、情報検索において使用される用語頻度重み及び逆ドキュメント頻度重みの汎用化である。それらは次のように定義される。

## 【数1】

$$tf_{ci} = \log(1 + N_{ci}) \text{ and } icf_i = \log \frac{N}{N_i}$$

ここで $N_{ci}$ は、コンテキスト $c$ 中の要素 $i$ の発生回数であり、 $N_i$ は、 $i$ が発生するコンテキストの数であり、 $N$ は、コンテキストの全数である。テキストモダリティ、ワードの対応する要素、及びドキュメントに対応するコンテキストの場合、この定義は、情報検索フィールドにおける用語頻度重み及び逆ドキュメント頻度重みに対する標準の定義と一致する。

40

## 【0054】

従って、テキストベクトルは、上述のようにトークン頻度重みを最初に計算し（ステップ410）、次に、上述のように、逆コンテキスト頻度重みを計算し（ステップ412）、次にテキストコンテンツベクトルを計算するためにこれら二つを乗算する（ステップ414）ことによって、計算される。

## 【0055】

50

本発明によって使用される埋め込みのためのトークン頻度重み及び逆コンテキスト重みの使用は、以下の直感的記述の一致する。コンテキスト（例えば、ドキュメント）中の要素（例えば、又はワード）の各異なる発生は、記述的特徴として、その要素に対する重要性の増加されたレベルを反映する。しかしながら、この増加は、線形である必要はないが、幾分"減衰"される。従来、対数が減衰関数として使用されたが、それはこの用途に対しても満足出来るものであることが判っている。同様に、逆コンテキスト頻度重みは、全てのコンテキスト（一例はテキストドキュメント中のワード"the"であるかもしれない）中に発生する要素に対する0から唯一つのコンテキスト中に発生する一つの要素に対する最大( $\log N$ )に達する範囲に亘る。対数スケールリングに対する一つのモチベーションは、情報理論に基づく。 $\log N / N_i$ は、コンテキスト中の要素*i*の発生について学習する時にいくらか多くの情報が得られるかの尺度として解釈され得る。ワード"the"が一つのドキュメントに発生することが学習される場合、（それが全てのドキュメントで発生すると仮定されると）重要な情報は得られない。しかしながら、フレーズ"Harry Truman"が一つのドキュメントで発生することが学習される場合、（そのフレーズが2, 3のドキュメントのみで発生すると仮定すれば）多くの情報が提供される。

【0056】

逆コンテキスト頻度重みが乗算されるトークン頻度重みは、ベクトルをスケールリングするための有利な方法であることがわかったことが理解されるべきである。しかしながら、他の重み付けスキームもまた可能であり、他の利点を提供し得る。

【0057】

従って、上述のように、テキストベクトルが計算されると、二つのテキストベクトル間の類似度が簡単なコサイン距離を介して計算され得る。

【数2】

$$\text{sim}_t(d_1, d_2) = \frac{\sum_i \phi_t(d_1)_i \phi_t(d_2)_i}{\sqrt{(\sum_i \phi_t(d_1)_i^2)(\sum_i \phi_t(d_2)_i^2)}}$$

ここで、 $d_1$ と $d_2$ は、二つの異なるドキュメントを表し、 $\phi_t(d_1)_i$ は、ドキュメント $d_1$ を表す*i*番目の項を表す。以下に更に詳細に議論されるように、ドキュメントの対同士の間のコサイン距離は、テキスト特徴単独に基づきドキュメントをクラスタリングするために使用されることが出来又は他の特徴と組合せて使用されることが出来る。

【0058】

本発明の他の実施の形態において、上述のテキスト特徴は、異なる方法で、或いは別の独立した特徴として、計算され得る。この他のバージョンにおいて、タイトル、ヘディング及びキャプションからのテキストのみが $R^{n_s}$ （ここで、 $n_s$ は、コレクション中のドキュメントのタイトル、ヘッダー及びキャプション中の特異のワードの全数である）中の"主題"のモダリティを定義するためにドキュメントから分離される。この他の（又は追加の）モダリティが上述されたテキストモダリティ（ドキュメントの完全なテキストのサブセットのみ殻のものを除き）と正確に同じ方法で引き出されるので、対応する特徴ベクトルと類似度を引出すために使用される公式は同じままである。

【数3】

$$\phi_s(d)_i = tf_{di}cf_i \text{ and } \text{sim}_s(d_1, d_2) = \frac{\sum_i \phi_s(d_1)_i \phi_s(d_2)_i}{\sqrt{(\sum_i \phi_s(d_1)_i^2)(\sum_i \phi_s(d_2)_i^2)}}$$

両実施の形態は、有用であることが検出され、必要ならば、交換的に又は共に使用され得

10

20

30

40

50

る。例えば、(例えば、タイトル中のワードの各発生を、それがそのテキスト中に2回又は3回発生したかのように、処理することによって)タイトル、ヘッディング及びキャプションテキストをドキュメント中の他のテキストとは異なるように重み付けすることも可能である。一般的な提案として、ドキュメント中の全てのテキストがテキストベースのモダリティの目的のために同じように処理される必要は無いことが認識されるべきである。調節と重み付けが可能であり、幾つかの用途で利点があり得る。

【0059】

同様に、ベクトルは、ドキュメントのURLに対して計算され得る。上述の例を詳論すると、その例示のURL "http://www.server.net/directory/file.html" は、7の用語 "http"、"www"、"server"、"net"、"directory"、"file"、及び "html" を含む。テキスト特徴に関して、これらの用語の幾つかは、殆ど又は全く情報的値を含まない(この例では、"http"、"www"、"net"及び"html")。従って、トークン頻度重み及び逆コンテキスト頻度重み埋め込みがここでは十分適切である。再度図4を参照のこと。

【0060】

この結果、各ドキュメント  $d$  は、 $R^{n_u}$  (各ベクトル空間は  $n_u$  次元を有し、ここで、各次元は、実数によって表される) に埋め込まれ、ここで、 $n_u$  は、コレクション中の全てのドキュメントを識別する特異のURLの全数を表す ( $n_u$  は、"URL要素の数"を表す)。ベクトル空間への埋め込みは、以下のように表される。

$$u(d)_i = tf_{di} icf_i$$

ここで、 $d$  は、特定のドキュメントであり、 $i$  は、ワードのインデックスであり、 $u(d)_i$  は、ベクトル  $u(d)$  の要素  $i$  である。トークン頻度重み ( $tf$ ) 及び逆コンテキスト頻度重み ( $icf$ ) は、情報検索で使用される用語頻度重み及び逆ドキュメント頻度重みである。それらは、以下のように定義される。

【数4】

$$tf_{ci} = \log(1 + N_{ci}) \text{ and } icf_i = \log \frac{N}{N_i}$$

ここで、 $N_{ci}$  は、コンテキスト  $c$  中の要素  $i$  の発生回数であり、 $N_i$  は、 $i$  が発生するコンテキストの数であり、 $N$  は、コンテキストの全数である。URLモダリティの場合、要素はURL用語に対応し、コンテキストは、ドキュメントに対応する。

【0061】

同様のベクトル埋め込みは、インリンクモダリティ ( $i(d)_i = tf_{di} icf_i$ ) 及びアウトリンクモダリティ ( $o(d)_i = tf_{oi} icf_i$ ) のために使用される。インリンクベクトルは、 $R^{n_i}$  に存在し、ここで、 $n_i$  は、コレクション中に埋め込まれた個別のインリンクの全数 (即ち、コレクション中のドキュメントの全数は、コレクション中の他のドキュメントに関連する)。アウトリンクベクトルは、 $R^{n_o}$  に存在し、そこで、 $n_o$  は、コレクション中に埋め込まれた個別のアウトリンクの全数 (即ち、コレクション中又は外において、ドキュメントの総数は、コレクション中のドキュメントによって参照される)。コサイン類似度は、類似的に計算される。

【数5】

$$\text{sim}_i(d_1, d_2) = \frac{\sum_i \phi_i(d_1)_i \phi_i(d_2)_i}{\sqrt{(\sum_i \phi_i(d_1)_i^2)(\sum_i \phi_i(d_2)_i^2)}} \text{ and } \text{sim}_o(d_1, d_2) = \frac{\sum_i \phi_o(d_1)_i \phi_o(d_2)_i}{\sqrt{(\sum_i \phi_o(d_1)_i^2)(\sum_i \phi_o(d_2)_i^2)}}$$

【0062】

本発明の他の実施の形態において、URL中の用語(上記のように定義されたURL埋め込みで使用されるように)インリンク及びアウトリンクから抽出され、そのように使用

10

20

30

40

50

された。しかしながら、この他の方法で引出された印リンクおよびアウトリンク特徴に基づくクラスタリングは、類似のドキュメントのクラスタリングにおいて効果は少ないことが判った。

【 0 0 6 3 】

ドキュメントのテキストジャンルは、 $R^{n_g}$ に埋め込まれ、ここで $n_g$ は、既知のテキストジャンルの数である。ドキュメントジャンルは、ドキュメントの解釈をガイドする文化的に定義されたドキュメントカテゴリである。ジャンルは、ドキュメントテキストではなくて、より大きなドキュメント環境（例えば、ニューヨークタイムズからナショナルエンクイアリを一目で区別するように働く物理的メディア、ピクチャー、等のような）によって合図される。二つの異なるジャンル中に存在する同じ情報は、二つの異なる解釈を導くかもしれない。例えば、ロー"At dawn the street was peaceful . . ."から始まるドキュメントは、小説の読者とは異なるようにタイムマガジンの読者により解釈されるであろう。各ドキュメントタイプは、容易に認識され且つ文化的に定義されるジャンル構造を有し、この構造は、それが含む我々の情報の理解と解釈をガイドする。例えば、ニュースレポート、新聞の社説、カレンダー、新聞発表、短編小説は、可能なジャンルの全ての例である。ドキュメントの構造とジャンルは、しばしばドキュメント又はテキストの自動解析によって（少なくとも一部）決定され得る（ステップ510）。テキストジャンルは、常に決定出来得るとは限らないが、特に、ウェブページ（それはしばしば良好に定義されたジャンルを有するわけではないが）の場合、多くの既知の可能なジャンルに対して、確率スコアのベクトルを計算することが一般的に可能であり（ステップ512）、次に、そのベクトルは、テキスト用語ベクトルに関して上述された方法で、類似度を決定するために使用され得る（コサイン類似度計算を介して）。

【 数 6 】

$$\text{sim}_g(d_1, d_2) = \frac{\sum_i \phi_g(d_1)_i \phi_g(d_2)_i}{\sqrt{(\sum_i \phi_g(d_1)_i^2)(\sum_i \phi_g(d_2)_i^2)}}$$

【 0 0 6 4 】

ベクトル空間に埋め込まれた画像に対して、二つのモダリティ（カラーヒストグラムと複雑さ）が良好に使用された。カラーヒストグラム特徴に対して、画像ドキュメントは、 $R^{n_h}$ に埋め込まれる。ここで $n_h$ は、ヒストグラム中の"ピン"の数（本実施の形態では12）である。好ましくは、一つの単一のカラーヒストグラムは、カラー特徴として使用される。この特徴空間は、HSV（その色相、彩度、及び明度のカラーモデル）へ変換され、二つのピンは、各次元に割り当てられる（ステップ610）。従って、カラー空間に対して三つの次元があり、各カラー次元に対して二つのピン（四つの値）があり、その結果、好ましいベクトル空間には、全12の次元がある。

【 0 0 6 5 】

次に、処理される画像中の各画素は、カテゴリに分けられ（ステップ612）、その色相、彩度及び明度は、各次元毎に4つのピンの一つに入り、対応するベクトル要素が増分される（ステップ614）。本発明の好適な実施の形態において、各ドキュメントに対するカラーヒストグラムは正規化され（ステップ616）、それによって、ピン値の全てが合計され、その結果がそのヒストグラムベクトルとして格納される（ステップ618）。テキスト（及び特定の他の）モダリティに対して行われるのが好ましいトークン頻度重みと逆コンテキスト頻度重み埋め込みを使用することは、このコンテキストでは意味がないので、適切ではない事を理解すべきである。しかしながら、ヒストグラムベクトル間の距離は、次のコサイン距離によって計算されることが好ましい。

【数7】

$$\text{sim}_h(d_1, d_2) = \frac{\sum_i \phi_h(d_1)_i \phi_h(d_2)_i}{\sqrt{(\sum_i \phi_h(d_1)_i^2)(\sum_i \phi_h(d_2)_i^2)}}$$

【0066】

本発明の他の実施の形態において、ヒストグラム間の距離は、最も大きなピン値による正規化で交点測定を介して計算され得る。

10

【数8】

$$\text{sim}_h(d_1, d_2) = 1.0 - \frac{\sum_i \min(\phi_h(d_1)_i, \phi_h(d_2)_i)}{\sum_i \max(\phi_h(d_1)_i, \phi_h(d_2)_i)}$$

【0067】

本発明の他の実施の形態において、複数のカラーヒストグラムは、各画像の複数の領域に対して決定され、複数のカラーヒストグラム特徴ベクトルを得る。例えば、4象限（左上部、右上部、左下部及び右下部）中のカラーヒストグラムと画像の中心は、別々に計算されることが出来、その結果、5つの別々のカラーヒストグラムベクトルを得、次に、ユーザによって必要ならば重み付けされ且つ組み合わされることが出来又は別々のベクトルとして残されても良い。或いは、上半分、下半分、左半分、右半分及び中心の矩形部分のような部分的に又は完全に重なる領域が使用され得る。効率化のためには、画像は、タイルに細分割され、ヒストグラムが各タイル毎に別々に計算され、次に、適切に領域へ組合せることが出来る。次に、画像の領域的類似度を介して画像を計算することが出来る。例えば、ブルースカイを有する全ての画像がそれらの"上部"カラーヒストグラムベクトルにおける類似度に基づき一緒にグループ化されてもよい。領域的画像類似度を求める他の実施の形態及び用途もまたここで記述される本発明のフレームワーク内で可能であることが理解されるべきである。

20

30

【0068】

これらの距離メトリックは、二つの画像に関して対称的である。対称距離は、一つの画像と他の一つの画像又は図心との間の距離が簡単な検索目的ではなく、クラスタリング目的のために必要なため、本フレームワークにおいて必要とされる。

【0069】

複雑さ特徴は、人間が画像同士間、一方では高いカラー均一性の領域を有する比較的少数のカラーから成る簡単なロゴと漫画と、他方では微細な陰影を有する比較的多数のカラーを有する写真との間、で行うことが出来る粗い意味的区別を捕獲するためである。この特徴は、画像内でカラーの水平及び垂直ラン長から引き出される。具体的には、同じカラーのラン（本実施の形態では、ステップ710において、上述のように、2ビットHSV値に粗く量子化される）は、x方向（ステップ712）及びy方向（ステップ714）に識別される。ヒストグラムは、各方向毎に計算され（ステップ716）、そこで、各ピンは、画素の数（他の実施の形態では、全高さ又は幅の量子化された割合）を表し、ランは、x方向及びy方向へ夫々スパンする。各ピンにおけるカウントは、その特定のラン長に属する画像中の画素数である。或いは、各ラン毎にピンへ追加される値は、そのランの長さによって重み付けされることが出来、より長いランへはより大きな重みを与えることが出来る。ヒストグラム中の要素の全数は、夫々画像の水平及び垂直次元における画素数である。従って、二つのベクトル（各ヒストグラムに対して一つ、水平方向及び垂直方向）が生成され（ステップ718と720）、画像複雑さに対する水平ベクトル及び垂直ベク

40

50

トルは、 $R^{n_x}$ （ここで、 $n_x$ は、画像の最大水平画素次元である）と $R^{n_y}$ （ここで、 $n_y$ は、画像の最大水平画素次元である）に夫々埋め込まれる。

【0070】

本発明のこの好適な実施の形態において、ラン長複雑さ情報は、少数のピン（従って、各ベクトル毎により少数の次元）に量子化される。これは、ベクトルの希薄性を減少し、画像間のより効率的及びより頑強な比較を可能とするために実行される。N個のピンと $n_x$ の最大水平次元が与えられると、 $n_x/4$ よりも長いあらゆる水平方向ランは、 $N^{1/h}$ （又は最後）ピンに配される。より短いラン $r_x$ は、フロアー $(r_x(N-1)/(n_x/4)+1)$ によってインデックス付けされたピンに配される（ここで、“フロアー”関数は、その項が最も近い整数に丸められる）。従って、ラン長は、N個のピンに線形に量子化され、 $n_x/4$ よりも大きな長さの全てのランは、最後のピンに入る。同様な操作が垂直方向ランに実行され、その結果、水平方向複雑さベクトルは、N次元を有し、垂直方向複雑さベクトルもまたN次元を有する。

10

【0071】

以下で述べられるように使用されるコサイン距離メトリックで、ピンの合計を正規化する必要はない。

【数9】

$$\text{sim}_c(d_1, d_2) = 0.5 \frac{\sum_i \phi_x(d_1)_i \phi_x(d_2)_i}{\sqrt{(\sum_i \phi_x(d_1)_i^2)(\sum_i \phi_x(d_2)_i^2)}} + 0.5 \frac{\sum_i \phi_y(d_1)_i \phi_y(d_2)_i}{\sqrt{(\sum_i \phi_y(d_1)_i^2)(\sum_i \phi_y(d_2)_i^2)}}$$

20

ここで、 $x$ と $y$ は、夫々水平方向複雑さベクトルと垂直方向複雑さベクトルを表す。

【0072】

或いは、二つのベクトル（水平方向と垂直方向）は、以下の標準コサイン距離メトリックを使用して、 $R^{n_x+n_y}$ （又は量子化された好適な実施の形態では、 $R^{2N}$ ）中のより大きなベクトルへ追加され得る。

【数10】

$$\text{sim}_c(d_1, d_2) = \frac{\sum_i \phi_c(d_1)_i \phi_c(d_2)_i}{\sqrt{(\sum_i \phi_c(d_1)_i^2)(\sum_i \phi_c(d_2)_i^2)}}$$

30

ここで、 $c$ は、追加されたベクトルを表す。

【0073】

カラー複雑さ特徴とカラーヒストグラム特徴の両方に対して、サブサンプリングがベクトル埋め込みの計算において生ずる演算費用を減少するために実行される。例えば、有用な結果を達成するように、画像中の端数（1/10のような）又は画素の全数の制限された数（1000のような）を選択することが可能であることが理解された。これらのサブサンプリングされた画素は、その画像全体に亘って均一に離間されることが好ましいが、他の実施の形態においては、ランダムに選択されてもよい。ヒストグラム特徴に対しては、サブサンプリングされた画素のみに対して適切なヒストグラムを計算するだけで十分である。複雑さ特徴に対しては、サブサンプリングされた画素が属する水平方向及び垂直方向のランの長さを決定することが必要である。本発明の好適な実施の形態において、これは、ロー及びコラムをサブサンプリングすることによって達成される。水平方向複雑さベクトルに対しては、最大50個の略等しく分布されたローの画素が選択され（画像が高さ方向において50画素よりも短いならば50未満）、且つこれらのローのみのランがカウントされる。同様の処理が垂直方向の複雑さベクトルのコラムに対して行われる。そうで

40

50

ないものは、ベクトル埋め込みが同じままである。

【0074】

最後に、（その他のベクトル埋め込みが指示されたように）ドキュメント間ではなくて、ユーザ母集団中のユーザ間の違いをハイライト可能な類似の特徴がある。例えば、ページ使用は、ユーザの情報探求優先度を指示することがわかった。ページ使用モダリティに対しては、ページアクセスは、最初に識別される（ステップ810）。好ましくは、トークン頻度重み（ステップ812）と逆コンテキスト頻度重み（ステップ814）が再び使用され、コンテキストは、各ユーザであり、且つトークンは、ユーザのページアクセスである。積はページ使用ベクトルとして格納される（ステップ816）。従って、ページ埋め込みは、 $\phi_p(u)_i = tf_{di}cf_i$ であり、ここで、 $u$ はユーザを表し、 $i$ はページを表す。その結果、埋め込みは、 $R^{n_p}$ になされ、ここで $n_p$ はコレクション中のドキュメントの全数である。他の実施の形態において、各ユーザのページアクセスは、2値とみなされる。それは、ユーザがページにアクセスした、その場合その対応するユーザのベクトルは対応する要素で"1"を有し、或いは、ユーザがページにアクセスしなかった、その場合対応する要素は"0"である。何れの場合においても、以下のコサイン距離メトリックは、（ユーザのページレファランスに関して）ユーザ間の類似度を計算するために使用され得る。

【数11】

$$\text{sim}_p(u_1, u_2) = \frac{\sum_i \phi_p(u_1)_i \phi_p(u_2)_i}{\sqrt{(\sum_i \phi_p(u_1)_i^2)(\sum_i \phi_p(u_2)_i^2)}}$$

【0075】

他のモダリティもまたユーザから導出され得る。例えば、ユーザ指定の人口統計情報（名前、年、ホビー、電話番号、ホームアドレス、選択されたグループのメンバーシップ等）及び他の種類のトラックされた情報（限定されるわけではないが、オンライン購買習慣、ソフトウェア使用及びドキュメントを見るために費やされた時間）もまた、スカラー又はベクトル空間に埋め込まれることが出来、（以下に述べられるように）数値距離メトリックの使用及びクラスタリングの実行が可能となる。例えば、ユーザのグループのメンバーシップは、既知のグループの数に等しい次元の数を有するベクトル空間に埋め込まれることが出来、ユーザのグループメンバーシップは、ユーザが対応するグループのメンバーであるか否かを表すブール（"0"又は"1"）を有する。これらの追加の例示のモダリティは、ここでは詳細には述べられない。しかしながら、本発明に係るシステムは、マッピングをベクトル空間に定義することによって、これらのモダリティや、ほぼあらゆるドキュメントベースの又はユーザベースの情報を組み入れるように容易に向上され得る。

【0076】

各モダリティ毎のベクトル空間中の次元の数は、ファクタの数によって変化し得ることを理解すべくである。例えば、テキストモダリティに対しては、各テキストベクトルは、コレクション中の特異のワードの数に等しい数の次元を有し、画像の複雑さモダリティに対しては、各ベクトルは、コレクション中の最大水平方向と垂直方向画素次元に等しい数の次元を有し、そして、ページ使用モダリティに対しては、各ベクトルは、コレクション中のドキュメントの数に等しい数の次元を有する。従って、ドキュメントがコレクションに追加すると（且つユーザがユーザ母集団へ追加すると）、同じ特徴に対するベクトルの全てが同じ次元を有し、それにより上述の類似度メトリックの使用を可能とすることを確実にするために、特徴ベクトルの多くを再計算することが必要となる。従って、演算費用を減少するために、ある状況では、有意な数のドキュメント（又はユーザ）が追加されるまで、特徴ベクトルのデータベースの更新を遅延することが有利である。勿論、新たなドキュメント（及びユーザ）は、それらが追加され、対応する特徴ベクトルが計算されるま

で、本発明に係るシステムによって認識されない。

【0077】

種々のモダリティの前述の表示は、本発明のシステムにおいてドキュメントとユーザ間の類似度をトラッキングするために有用であり効率的であることがわかった。しかしながら、ドキュメント情報のベクトル空間への埋め込み及びドキュメント間の類似度の計算方法の種々の他の方法も可能であることが認識されるべきである。例えば、ドキュメントに対応するテキスト、URL、インリンクテキスト及びアウトリンクテキストを単一の支配的なテキストベクトルへ組合すことが可能である。このアプローチは、画像ドキュメントと関連するほんの僅かなテキストがある時、有用であり得る。また、上述のコサイン類似度メトリックは、一度に、単一の特徴又はモダリティに基づいて、二つのドキュメント間の類似度を計算することが理解されるべきである。また、二つのドキュメント間の集合類似度を計算することが、可能であり、状況下で好ましい。

10

【数12】

$$\text{sim}(d_1, d_2) = \sum_j w_j \text{sim}_j(d_1, d_2)$$

ここで、jは上述の適用可能モダリティを表すと共にその範囲に亘り、 $w_j$ は各モダリティに対応する重み付けファクタ（好ましくは、単一ではあるが、必要に応じて調節可能である）を表す。次に、この集合類似度は、全ての可能な（又は実際の）モダリティの基づくドキュメント間の全体の類似度を表す。

20

【0078】

全てのモダリティが全てのドキュメント中にあるわけではないことが前述より明らかである。例えば、ウェブ（又は、ウェブのようなイントラネットコレクション）上で、テキスト、画像又は全体として他の何かである、全てのドキュメントは、検索のためにドキュメントを識別するように働く一つの対応するURLを有する。しかしながら、ドキュメントの全てが画像ではなく、全てのドキュメントが画像ではなく、それにより、ヒストグラムと複雑さメトリックが幾つかのドキュメントに対して不能である。同様に、全てのドキュメントがテキストを含むわけが無いが、（上述のように）テキストは、ある場合（インリンクがある場合）ドキュメント参照から組み立てられる。

30

【0079】

従って、集合類似度メトリックは、ある状況で、サブ最適であってもよく、必要な時、個々の類似度メトリック上に"フォールバックする"能力を有することが望ましい。  
クラスタリング

【0080】

集合類似度メトリックを含む、上述の類似度メトリックは、ドキュメントとユーザ（集合的にオブジェクト）をクラスタリングするためのベースを定義する。標準のクラスタリングアルゴリズムが使用される。本発明のこの好適な実施の形態において、"k平均"クラスタリングがk個の異なるクラスタへオブジェクトを割り当てるために使用される。

【0081】

40

当該技術で公知のように、k平均クラスタリングは、通常クラスタセンターとしてk個のランダムに選択された複数のオブジェクトで始まる区分方法である。複数のオブジェクトは、クラスタセンター（オブジェクトが最も高い類似度を有するセンター）へ割り当てられる。次に、クラスタセンターはそれらのメンバーの平均として再計算される。オブジェクトの（再）割り当てと平均の再計算の処理は、それが収束するまで数回繰り返される。クラスタの数kはその方法のパラメータである。k=20及びk=50の値は、それらの値が良好な結果を与えたので、種々の実施及び研究に使用されたが、との値もまたユーザの優先度に基づき等しい効果を得るために使用され得る。

【0082】

本発明の他の実施の形態において、階層マルチモードクラスタリングもまた使用できる

50

が、k - 平均クラスタリングが満足な結果を提供することが確認された。

【0083】

上述のように、k 平均クラスタリングの旧来のフォームは、クラスタリングされるべきオブジェクトからのランダム選択によって初期クラスタを選択する。初期クラスタを選択するための他の方法は、バックショットアルゴリズムを使用し、階層（しかし演算的に効果である）クラスタリングアルゴリズムをそれらのオブジェクトのサブセットに適用することによって初期センターを計算する。次に、k 平均クラスタリングに対する初期センターは、そのサブセットをクラスタリングすることにより発見されるクラスタのセンターである。

【0084】

しかしながら、ランダム選択及び階層サブセットクラスタリングの両方は、マルチモードクラスタリングに対してはサブ最適であることが判った。ドキュメントコレクションでは典型的なベクトル空間は、しばしばその空間の一つの小さな領域に共に集中する大部分のオブジェクトと、他の領域に希薄に分布する他の有意な数オブジェクトとを有する。このタイプのデータに対しては、初期センターを識別するためにウエーブフロントクラスタリングが一層効率的であることが判った。このウエーブフロントアルゴリズムは、以下のように且つ図9に示されるように進む。

【0085】

最初に、m個（クラスタリングされるべきオブジェクトの総数Nよりもかなり小さい数）オブジェクトがランダムに選択される（ステップ910）。この数は、数k（それは、最終的に計算されるクラスタの数である）から独立している。実験により、mの適切な値は、10であることが判った。

【0086】

次に、m個のオブジェクトのベクトル図心aaaが計算される（ステップ912）。この図心は、当該技術で公知の方法によって、即ち、対象となるベクトルの対応する項を平均することによって、計算される。

【0087】

次に、トータルでk個のオブジェクトbbbがクラスタリングされるべきN個のオブジェクトからランダムに選択される（ステップ914）。上述のように、kは、最終クラスタの望ましい数である。最後に、k個の初期オブジェクトbbbの各々に対して、k個の初期オブジェクトの各々への図心aaa回りにk個のクラスタセンターcccを計算する。これらのクラスタセンターは以下のように計算される（ステップ916）。

【数13】

$$\bar{x}'_i = \alpha \bar{c} + (1 - \alpha) \bar{x}_i$$

ここで、 $i = 1 \dots k$  に対して、 $\alpha$  の適切な値は0.9であることが判った。他の値でも有効である。

【0088】

この技術は、名称"ウエーブフロントクラスタリング"が与えられた。理由は、単純化された項において、"ウエーブ"は、図心aaaから送られ、且つこのウエーブによって第2のセットのランダムにピックアップされたオブジェクトへの途中でヒットされるオブジェクトが初期クラスタセンターとして選択されるからである。これらの初期センターは、多数のオブジェクトが一つの点に収束する場合に適切である。理由は、図心aaaは、その点へ近接している傾向があるからである。これらの初期センターは、集中する領域を効率的に区分するのに良好に当てはまる。

【0089】

次に、標準のk - 平均クラスタリングは、図10に示されるように、各オブジェクトをその最も近いクラスタへ割り当てられることによって進む。最初に、図9に図示されているようにクラスタセンターを選択した（ステップ1010）後、割り当てられなかったオ

10

20

30

40

50

プロジェクトが選択される(ステップ1012)。その類似度は、上述の類似度メトリックの一つを使用して、各クラスタセンターに関して計算される(ステップ1014)。次に、そのオブジェクトは、最も近いクラスタセンターへ割り当てられる(1016)。割り当てのためのオブジェクトがもっとある場合、その処理を繰り返す(ステップ1018)。次に、クラスタセンターは、各クラスタセンターに対応する各クラスタの図心(平均)として再計算される(ステップ1020)。クラスタセンターが、例えば、十分に少数のオブジェクトがクラスタにスイッチされたか否かを決定することによって、十分に収束した(ステップ1022)時、クラスタリング処理が終了する(ステップ1024)。或いは、全てのオブジェクトは、すべてのクラスタから割り当てが解除され(ステップ1026)、処理は、再び新たに決定されたクラスタセンターで始まる。

10

【0090】

用途

本発明のシステム及び方法を示すため、マルチモード特徴の二つの用途がここで考察される。(1)マルチモードブラウジング及び検索と呼ばれるシステムにおいて対象となるドキュメントをユーザが識別することを助けることと(2)ユーザのコレクションとの対話のマルチモード解析(コレクション使用解析、即ちCUA)である。

【0091】

第1の用途において、上述されたように生成されたドキュメントのクラスタは、ドキュメントのサーチ、推奨及びブラウジングのためのシステムにおいて使用される。第1の用途の第1の実施の形態において、一つの特徴は、ユーザの指定のように、一度に考えられ、第2の実施の形態において、複数の特徴が同時に考えられる。

20

【0092】

第2の用途において、上述のように生成されたユーザクラスタは、二つの別々の関数に適用される。第1に、ユーザクラスタは、調停を介して可視化に適するようにされ、これについては、以下でより詳細に記述される。第2に、マルチモードユーザクラスタが推奨を生成するために使用される。

【0093】

以下、これら二つの用途において、マルチモード情報の使用が記述される。その方法は、このような情報を組合せ且つ例を介してそれらの利点を示す方法を含む。

【0094】

シーケンシャルマルチモードブラウジング

一度に一つのタイプの特徴を使用してのマルチモードサーチとブラウジングは、図11乃至22に関連して最も良く示される。各特徴は、そのセットの画像をリファインするか対象となる画像の関連セットへマッピングするために使用される。従って、画像特徴は、テキスト特徴から独立して使用され、テキストがより適切に知覚される時のテキスト(例えば、セクションヘッディング、アブストラクトタイトル、画像アンカー中の"ALT"タグ)の使用と、画像特徴がそのようにより多い時の画像特徴の使用との間を人間のユーザがナビゲート出来る複数のクラスタリングを生成する。

30

【0095】

異なる特徴に基づきサーチを徐々に狭めることに一つの起こり得る問題は、特徴値が欠落した画像が考慮から不注意で削除されるかもしれないことである。例えば、幾つかのドキュメントは、テキストに関係無い画像や、画像のコンテンツに関連しないテキストを画像を含む。具体的には、幾つかの画像は、テキストの無いページ上に存在する。他の場合、画像を囲むテキストはその画像の意味的コンテンツに関連しない。サーチを徐々に狭めることの他の問題は、そのサーチが二つのクラスタ間の境界近くの空間の一部へ狭められるかもしれないことである。

40

【0096】

特徴の使用は、ここで対象となるセットの要素への迅速な初期フォーカシング、次の編成及び類似の要素を含むための拡張を可能とし、類似の要素の幾つかは、不完全な特徴セットを含んでも良く、或いは他の一つのクラスタに生じてもよい。

50

## 【0097】

ここで提案される方法の幾つかは、画像ブラウジングへの拡張として考えられ得る。理想的なブラウジングシステムは、ユーザが画像を含むドキュメントのブラウジングを可能とし、それらのドキュメントは、記述的注釈テキストを含んでも含まなくてもよく、且つテキスト特徴及び又は画像特徴を使用する。ユーザは、意味的コンテンツ("その画像が何を示すか?")か視覚的コンテンツ("その画像がどのように見えるか?")に基づいて画像コレクションがブラウジングすることを望むことが出来る。画像検索システムは、画像を意味的情報で自動的に注釈することは、現在不可能なタスクであるために、しばしば手作業によるキーワード注釈や画像特徴のマッチングに基づく。そうとしても、手作業によりラベル付けされた画像コレクションは、画像が有することが出来る全ての可能な意味的重要度を含むことは出来ない。

10

## 【0098】

上述のように、ここで記述されたアプローチは、Cutting等の論文で述べられたScatter/Gather方法といくつかの点で類似する。

## 【0099】

Scatter/Gatherは、元来はドキュメントから引出されるテキスト特徴と共に使用されるようにデザインされている。Scatter/Gatherは、コレクションを少数のクラスタに"分散する"ことによってサーチを繰り返しリファインし、次に、ユーザは、再び分散するために対象となるクラスタを"集める"。Scatter/Gather方法は、本発明によってテキスト特徴と画像特徴の両方を使用してテキストと画像を有するドキュメントのコレクションをナビゲートして、拡張される。また、ワーキングセットの外側からの要素がワーキングセットに組み込まれることが出来るように"拡張"(即ち、マッピング)機能がある。

20

## 【0100】

マルチモードブラウジング、推奨及び可視化への本発明のアプローチにおいて、照会に対する正確な解答は、ユーザに依存する。従って、ブラウジングに関連する本発明の一態様において、ユーザは、各ステップで使用される特徴を選択する。ユーザは、現在のワーキングセットを見るのみである。マップ機能が使用されず、各オペレーション後に唯一のクラスタが選択される時、これは、ユーザがデプスファーストサーチ中のツリーの唯一のノードを拡張することに等しい。クラスタを組合せるために選択することによって、格子が形成される。そして、マップ機能を使用することによって、ワーキングセットの外側からの要素は、ワーキングセットの一部となってもよく、それによって、ツリーも格子も生成されない。従って、本発明の方法は、決定ツリーとは全く異なっている。

30

## 【0101】

実際、初期テキスト照会は、対象となる候補画像を検出するために使用され得る。次に、対象となる画像を含む戻されたクラスタの幾つかは、更なる考察のためにユーザによって識別される。一つの画像特徴の類似度に基づく拡張によって、システムは、関連するテキストなしに又はユーザ指定の照会に対して十分には類似しないテキストと共に、初期に選択されたクラスタによって表示されるものに類似する画像クラスタを検出及び提示する。このように、拡張機能は、テキスト照会の結果として、オリジナルセット中にない等価画像が識別され且つ含まれることを許容する。拡張機能は、また対象となる特徴空間に近いが、前のステップでの区分に起因して他の一つのクラスタ中にある要素を考察のために識別できる。

40

## 【0102】

上述のように、本発明のマルチモードブラウジングと検索態様において、事前処理ステップがブラウジングの間に必要な情報を事前演算を行い、そのデータの初期編成を提供するために使用される。一セットの個々の特徴(多分、異なるモダリティからの)は、各ドキュメント毎に事前演算され、ベクトルとして格納される。本出願において、ウェブページ中の画像の特徴は、以下のように演算される。テキスト特徴は、各画像を囲みそれと関連するテキスト、その画像のURL, ALTタグ、ハイパーリンクテキスト、及びテキス

50

トジャンル（後述される）を含む。画像特徴は、カラーヒストグラムとカラー複雑さのメトリックを含む。上記の表 1 を参照。ドキュメントは、特徴の各々に基づいてグループにクラスタリングされる。

#### 【 0 1 0 3 】

画像のサーチのために、ユーザはテキスト照会を入力することで開始する。仮説セッションが図 1 1 に示されており、その図で、円形ノードは、クラスタ中のデータを表し、実線の矢印はノード中のデータの拡散（スキューリング）又は集合を表し、鎖線は、拡張（マップ）機能におけるように、一つのノード中のサブセットのデータの他のノードへの移動を表わす。事前演算されたテキストクラスタは、コサイン距離を使用して照会用語への等価（即ち、類似度）に関してランク付けされ、最高にランク付けされるクラスタが戻される。これらは、第 1 のセットの結果 1 1 1 0 に代表的テキスト又は画像としてディスプレイされ得る。次に、ユーザは、ユーザの対象に最も類似するクラスタを選択する。これは、全て又はサブセットのクラスタ 1 1 1 2 を含み得る。次に、以下の二つの操作の一方が一般的に実行される。選択されたクラスタ中の画像は、選択された特徴に基づき再クラスタリングが行われ、結果として他のセットの結果 1 1 1 4 を得る。又は選択されたクラスタが選択された特徴に基づき新たな類似クラスタ 1 1 1 6 へマッピング（拡張）される。

10

#### 【 0 1 0 4 】

何時でも、ユーザは、新たなサーチを開始でき、又は新たな照会（初期テキスト照会のような）を実行することによって既存の結果セット上に動作することが出来る。次に、後者の照会の結果は、ユーザのオプションとして、その既存の結果セットをリファイン又はそれへ追加するために使用され得る。

20

#### 【 0 1 0 5 】

新たなクラスタは、選択された特徴がテキストデータから引出されたか又は画像データから引出されたかに依存して、代表的テキスト又は画像としてディスプレイされる。選択された特徴は、事前演算された特徴の何れであってもよい。再クラスタリングすることによって、ユーザは、そのセットの画像をリファインすることが出来る。マッピング又は拡張（即ち、先行するリファインが行われるか否かに拘らず同じ又は類似クラスタにおいて他の類似のドキュメントを識別すること）によって、指定された特徴において類似するが多分他の特徴において値が欠落している画像考察のためにそのセットの画像へ持ち込むことが出来る。

30

#### 【 0 1 0 6 】

上述のように、クラスタリングは、予め設定された数のクラスタを有する標準  $k$  - 平均クラスタリングを使用して実行される。上述の事前演算ステップにおいて、クラスタの数は、ユーザに提示されるクラスタの数よりも大きい。これは、サブセットのクラスタのみが初期テキストストリング照会に回答して提示されるためである。初期テキスト照会を有する本発明の一実施の形態において、12個のクラスタは、最初に使用されるが、5個の最も類似するクラスタのみがその照会に基づいて戻される。次に、集めるためにユーザによって選択されたクラスタは、再クラスタリングされ、そこで、クラスタの数は、ディスプレイされるべきクラスタの数、本実施の形態では、5個、に等しい。各異なる集合とクラスタリング操作は、5個のクラスタを得る。各操作が実行されると、クラスタの結果は格納される。これは、一連の操作をバックアップすることを可能とし、それはまたマッピング又は拡張操作によって必要とされる。

40

#### 【 0 1 0 7 】

初期クラスタリングは、或いはカラーヒストグラム特徴のような一つの特徴に基づいてもよい。初期クラスタの適切な数は、特徴に依存してより少なくても良い。この実施の形態において、初期クラスタリングは、テキストに基づくが、何時でも、拡散と更なるクラスタリングは、テキスト特徴又は画像特徴に基づいてもよい。本発明の他の実施の形態において、非テキスト特徴に基づく初期クラスタリングが可能であり、且つそれはある状況において有用であり得る。

50

## 【 0 1 0 8 】

上述のように、拡張/マップ機能は、異なる特徴に基づくサーチを徐々に狭めることの問題を扱い、ここにおいて、値が欠落した画像が考察から削除される。例えば、幾つかのドキュメントは、関連するテキストを有さない画像やその画像のコンテンツに関連しないテキストを含む画像を含む。他のケースでは、画像を囲むテキストは、画像の意味的コンテンツと関連しない。サーチを徐々に狭めることの問題は、そのサーチが二つのクラスタ間の境界に近い空間の一部へ狭められてもよいことである。

## 【 0 1 0 9 】

マッピング又は拡張操作は、画像又はクラスタを一つの特徴次元における類似度に基づき現在のセットへ追加する。唯一の特徴が一度に考慮されるので、類似度確立に使用される距離メトリックが各特徴毎に異なっても良い事を理解すべきである。例えば、上述のように、コサイン距離は、テキスト特徴類似度のために使用されることが出来、ユークリッド距離又は正規化されたヒストグラム交差がヒストグラム類似度のために使用される。

10

## 【 0 1 1 0 】

拡張操作は、幾つかの方法で実行され得る。一つの方法は、現在のクラスタの要素がマッピングされたセット中に残り、そのセットサイズが増加されることを確実にする。これは、選択された特徴に基づき（適切な距離メトリックを介して）現在のワーキングセットへ近い幾つかの要素をそのワーキングセットへ追加することによって達成される。この好適な実施の形態において、現在のワーキングセットに対して選択された特徴の平均が計算され、次に、この平均に最も近い、全体のデータベースから選択されたこれらの要素（ベクトルとして表される）が追加される。これはテキスト特徴に対して最も適切である。他のバージョンでは、ワーキングセット中の各表示された代表に近い要素が選択され追加される。この他のマッピング手順は、画像特徴により多く適用可能であり、そこで、クラスタは、テキストを表すために使用される要素の収集物の代わりに選択された画像によって表示される。しかしながら、テキストが選択されたドキュメントによって表示される時、マッピングの後の方法は、また適切である。

20

## 【 0 1 1 1 】

マッピングは、上述されたように、バックアップのためにセーブされた連続したワーキングセットの更に上に提示されたこれらの要素のみを考慮することによってスピードアップされ得る。即ち、マッピングのために選択された特徴がクラスタリングのために使用されるまで、バックアップの連続する操作を参照する。例えば、クラスタリングがカラーヒストグラム特徴に基づき実行され、URL特徴に基づき更なるクラスタリングが続いたと仮定する。カラー複雑さに基づくマップ操作が必要な場合、カラーヒストグラム（他の一つの画像特徴）に基づき選択されたクラスタからの要素が、全てのクラスタではなくて、使用され得る。

30

## 【 0 1 1 2 】

最終の拡張は、特徴のためのデータを有さない要素の全てを含む各特徴毎の特別のクラスタを生成することを含む。マッピングが実行されるべき時、すでに使用された特徴を関連する特別のクラスタ中のこれらの要素のみが候補と考えられ、現在のワーキングセットへ追加される。

40

## 【 0 1 1 3 】

図 1 1 及び上述のカラーヒストグラム/URL特徴の例を参照して、マッピングのための他の一つの（より簡単な）方法は、カラーヒストグラム特徴に基づき最も類似するクラスタを識別することを含む。この方法では、関連するテキストを有さない画像は、それらが適切に関連するテキストを有する画像と類似する場合は、同一である。例えば、あるURLが情報を有さない（例えば、"`http://www.company.com/products/special/image.jpg`"、それは、唯共通の用語"`www`"、"`company`"、"`com`"、"`products`"、"`special`"、"`image`"、及び"`jpg`"を含む）。最初にURL特徴を有する画像を識別し、次に、他の一つの特徴において類似する画像へマッピングすることによって、より多くの数の画像がサーチ

50

を再開したり、特徴重みの使用を必要とすることなく識別され得る。

【0114】

Scatter / Gatherのようなクラスタリングスキームを使用する時、ブラウジングセッションの間ユーザに対してクラスタをディスプレイ或いは表示することが必要である。テキストクラスタは、多くの方法で表示されることが出来、最も一般的には、ある意味で、クラスタの最も代表的であるワードのセットの選択とディスプレイである。画像クラスタが表示されることが必要な時、クラスタメンバーに共通する画像特徴を選択しそれらをディスプレイすることは殆ど意味がない。理由は、これらは、一般的に、ユーザにとって意味的な意義を持たないからである。従来クラスタリング画像ブラウザは、画像をより低い(2)次元空間へマッピングしそのマップをディスプレイすることによって画像クラスタを表示していた。代わりに、本発明の好適な実施の形態は、クラスタの更なるクラスタリングを呼び出し、次に(a)そのクラスタの図心に最も近い三つの画像と(b)そのクラスタのサブ領域を代表する三つの画像によってそのクラスタを表示する。その三つのサブ領域の代表は、その三つの最も中心の画像を上記(a)から除去し、三つのサブクラスタを計算し、そして各サブクラスタの図心に最も近い画像(適切な距離メトリックを介して測定される)を使用することによって計算される。この表示は、クラスタ図心の向きとクラスタ中の画像の範囲を提供する。また、代表的画像は、多次元スケーリングを使用して2次元ディスプレイ上に配されてもよいが、本発明の例では、それらの代表は、一コラムの三つの"図心"画像又は三つの"サブクラスタ"画像にディスプレイされる(例えば、図14参照)。これにより、サムネル画像及びオリジナルの複数のコピーのような非常に類似する画像がより容易に識別され得る。

10

20

【0115】

2310個の画像を含むウェブのようなドキュメントのコレクションが、以下の例に対する例示のコーパスとして使用された。ウェブドキュメントは、ドキュメントの走査画像中に発見されることが出来、ドキュメントのコンテンツやドキュメント中の構成要素を推理するために使用出来る同じタイプの"メタ情報"の多くを含む。ウェブドキュメントと共に働くことによって、画像中の構成要素及びレイアウトを識別するのに伴う問題が最小化され、且つ検索処理でメタデータを使用するための技術の開発が可能となる。

【0116】

ウェブ上に遍在するロゴとアイコンのような"興味のない"画像によってコーパスが支配されることを防止するために、画像が満足しなければならない幾つかの単純で幾分任意の基準がコーパスに含まれるように適用された。何れかの特定のクラスの全ての画像を含むことが必要でも実行される実験のゴールでもなく、ウェブ上で利用出来る興味あるコーパスをアSEMBLすることのみが必要で、それによって高い拒絶閾値が意図的に使用された。画像は、少なくとも50画素の高さと幅を有することが必要であり、少なくとも10000個のトータル画素を含む事が必要であった。また、画像は、次の幾つかのカラーコンテンツベースのテストをパスすることが必要であった。その画像の90%以下が8カラーから成り、その画像の95%以下が16カラーから成り、画像の画素のRGBカラー空間共分散マトリックスが非単一であった。質的に、これらの基準は、その画像が単純なライン図ではなく、詳細に上述されたように、カラー特徴によって十分に区別出来得るような十分に種々のカラーコンテンツを含むことを保証する。同じ画像の複数のバージョンに対してスクリーニングが実行されなかった、それによってコーパスは同一の画像と、一つの画像とそのサムネル画像を含まない。

30

40

【0117】

異なるモダリティにおいて"分散"と"集合"の使用を示す三つのサンプルセッションが以下に述べられる。第1の例は、最初にコレクションを狭めるためにテキスト特徴を使用し次に結果を編成するために画像特徴を使用することを示す。最初に、図12を参照して、ユーザがテキスト照会"ancient cathedral(昔の大聖堂)"1210をタイピングし、"実行"ボタン1212を押すことによって開始する。ここで述べられるようにシステムとのユーザの対話が公知の方法で、例えば、実際の物理的ボタンとの相互作用

50

用により、マウスのようなポインティングデバイスと有するボタンのスクリーン上の表示の操作により、あまり多くはないがネームへのボイスコマンドより起こる。本発明のこの好適な実施の形態では、ユーザは、本発明を実施するソフトウェアプログラムによってウインドウ 1 2 1 4 として提供されるマルチモード画像ブラウザと対話する。

【 0 1 1 8 】

5つの戻されたテキストクラスタ 1 2 1 6、1 2 1 8、1 2 2 0、1 2 2 2 及び 1 2 2 4 をディスプレイするスクリーンのスナップショットが図 1 2 の左半分に示されている。これらのクラスタは、照会用語に最も近いクラスタである。各クラスタにおける最頻度コンテンツ用語は、各クラスタを表示するようにディスプレイされる。ユーザは、テキストクラスタに対する追加の表示用語を見るために各テキストウインドウをスクロール出来る。ユーザは、テキストに基づいて、用語 "a c i e n t (昔の)" と "c a t h e d r a l (大聖堂)" を含む第 1 のテキストクラスタを再び分散することを決定する。そうするために、ユーザは、望ましいクラスタの次のチェックボックス 1 2 2 6 を選択し、引き続いて "テキストクラスタ" ボタン 1 2 2 8 を押す。上述のように、これによって、そのシステムが既存の選択されたクラスタをより小さな別々のクラスタにリファインする。

10

【 0 1 1 9 】

5つの結果のテキストクラスタ 1 3 1 0、1 3 1 2、1 3 1 4、1 3 1 6 及び 1 3 1 8 をディスプレイするスクリーンのスナップショットが図 1 3 の左半分に示されている。ユーザは、用語 "a c i e n t"、"c a t h e d r a l" と "c h u r c h (教会)" を含む 3 つのクラスタを選択して、(対応するチェックボックス 1 3 2 0、1 3 2 2 及び 1 3 2 4 を介して) 集め且つ ("複雑さクラスタ" ボタン 1 3 2 6 を押すことによって) 分散のための特徴として複雑さを選択する。

20

【 0 1 2 0 】

画像の複雑さに基づくクラスタリング後のスクリーンのスナップショットが図 1 4 に示されている。図心最も近い表示画像がディスプレイされる。各画像クラスタの次の矢印 (例えば、第 1 の画像クラスタ 1 4 1 4 に対応する左矢印 1 4 1 0 と右矢印 1 4 1 2) をクリックすることによって、ユーザは、図心とサブクラスタ表示ビューとの間を移動出来る。画像クラスタ 1 4 1 4、1 4 1 6 及び 1 4 2 0 は、古い教会及び大聖堂を含む、"a n c i e n t" ビルディングとモニュメントを主とする画像を含む。画像クラスタは、1 4 1 8 は、ロゴを含み、画像クラスタ 1 4 2 2 は、種々雑多なアイテムを含むように現れる。

30

【 0 1 2 1 】

第 2 の例において、我々の仮説ユーザは、我々のコーパス中に多くの p a p e r m o n e y (紙幣) の画像を見つけ出すように試みている。図 1 5 に示されるように、"p a p e r m o n e y" の初期照会が与えられ、且つ結果としてのテキストクラスタ 1 5 1 0、1 5 1 2、1 5 1 4、1 5 1 6 及び 1 5 1 8 がディスプレイされる。第 1 のテキストクラスタ 1 5 1 0 は、ワード "m o n e y (お金)" とワード "n o t e (紙幣)" を含む。第 2 のテキストクラスタ 1 5 1 2 は、ワード "p a p e r (紙)" を含むが、それを囲むワードは望ましい意味のワード p a p e r が使用されていることを示さず、それによってこのクラスタは選択されない。お金が多くのカラでプリントされるので、最初に画像特徴としてカラー複雑さメトリックが使用されるのが適切である。従って、第 1 のテキストクラスタ 1 5 1 0 は、カラー複雑さ特徴に基づいて、第 1 のテキストクラスタ 1 5 1 0 が分散され、結果としてのクラスタが図 1 6 に示されている。画像クラスタ 1 6 1 4 と 1 6 1 8 は、紙幣の画像を含み、それによってそれらは、(両方のクラスタを選択することによって) 集められ、次にこの時にカラーヒストグラムに基づいて分散される。他の画像クラスタ 1 6 1 0、1 6 1 2 及び 1 6 1 6 は、対象となる画像を含むようには現れず、それによってユーザはこれらを選択しない。

40

【 0 1 2 2 】

結果の画像クラスタは図 1 7 に示されている。画像クラスタ 1 7 1 2 は、1 4 個の画像を含み、且つ中央の代表例は、紙幣の全ての画像である。このクラスタは、ヒストグラム

50

特徴に基づき再び分散される。図18に示されているように、それが紙幣の多くの画像を含むことが観察され得る。画像の幾つかは、複製であるように現れるが、この場合、それらは実際には一つのサムネル画像と一つのフルサイズ画像である。サブクラスタ表示の審査は、マネー（お金）を含まないサブクラスタ中の幾つかの画像を表すが、それらは、マネー画像と類似するカラーを有する。

【0123】

この例は、選択的にそのセットの画像を一セットの対象物に狭めるための一連のコンビネーションにおける異なる特徴の使用を示す。分散は、より大きなコレクションをより小さなサブセットへ編成することを助けるために使用される。集合することによって、異なるコレクションが組み合わされ且つ一緒に再編成される得る。

10

【0124】

最後の例において、図19に初めに示されているように、ユーザは、照会"pyramid egypt"（ピラミッドエジプト）においてピラミッドとタイプをサーチしている。戻されたテキストクラスタ1910、1912、1914、1916及び1918がディスプレイされる。ユーザは、複雑さ特徴に基づき分散されるべき第1のテキストクラスタ1910を選択し、結果としての画像クラスタからの代表画像が図20に示されている。ユーザは、第2及び第4の画像クラスタ2012と2016に石を含むアウトドアシーンがあることに気付く、カラーヒストグラム特徴に基づき、更なるクラスタリングのためにこれらを選択する。結果としての画像クラスタが図21に示されている。第1の画像クラスタ2110は、4つの画像を含み、且つ第1の画像はピラミッドの画像である。

20

【0125】

第1の画像クラスタ2110が、（第1の画像クラスタ2110を選択し、且つ"ヒストグラム拡張"ボタン2120を押すことによって）カラーに基づいて類似画像を含むように拡張されると、図22に示されているように、ピラミッドの他の画像2210が識別される。この画像は、テキスト無しで且つ意味のないURLと共にウェブページ上に発生し、それによって、カラーヒストグラム特徴に基づいてそれが検索された。

【0126】

この例において、テキスト照会は、画像コレクションのサイズを減少するために使用され、且つ減少されたコレクションは、画像複雑さ特徴に基づきプレゼンテーションのために編成された。カラーヒストグラム特徴次元で類似する追加の画像が得られた。

30

【0127】

これらの例において、異なるモダリティにおける特徴は、コーパス中のサブセットの要素の"分散"と"集合"の技術を使用して、関連するテキストを有する一セットの画像をユーザがブラウジングするのを連続して助けるために使用される。セッションは、テキスト紹介で始まり、全体のコーパスよりも多くの焦点を当てられた初期セットで開始する。一つ又はそれを越える興味ある要素を含むことが観察されるクラスタは、それらのコンテンツを審査するために分散されるか、又は全体のコレクションから類似の結果を検索するために拡張される。上述の例（図12乃至図22）は、三つの特徴タイプ、即ち、テキスト、画像ヒストグラム及び画像複雑さ、のみを使用したが、本発明の方法は、ここで記述される全8つのモダリティ、及び他にも等しく適用出来る。

40

【0128】

従って、本発明の一態様は、複数のモダリティを利用してコレクションをブラウジングするためのシステムを含む。クラスタを審査するために要素を"集合"且つ"分散"する繰り返し処理を介して、ユーザが対象となる画像のグループを検出出来る。"拡張"又は"マップ"機能は、一つ又はそれより多くの次元において値を欠落しているが対象となるある次元における他の要素に類似するコレクション中の要素の識別することを可能とする。

【0129】

集合マルチモードブラウジングまた、上で提案されたように、クラスタリングと拡張操作のために距離メトリックの種々の組合せを使用することが可能である。

【0130】

50

上述された例示のシステムと方法を使用してこれを実施するために、二つのドキュメント又はオブジェクト間の集合類似度  $sim(d_1, d_2)$  が、前述のセッションで記述された集合、分散及び拡張操作で使用され得る。図 1 2 乃至 2 2 に示されるユーザインターフェースへの僅かな変更は、この追加の特徴を許容する。例えば、"集合クラスタ"と"集合拡張"ボタンが、全ての可能なモダリティを同時に操作することを容易にするために追加出来、或いは可能なモダリティのリストアップ(テキスト、カラー複雑さ、カラーヒストグラム等)は、"クラスタ選択モダリティ"又は"拡張選択モダリティ"ボタンが起動されると、一つのモダリティが使用されるべきか、一度に複数のモダリティが使用されるべきかをユーザが指示することが出来るようにチェックボックス(及び任意ではあるが、ユーザ調節可能重み)が設けられ得る。次に、選択されたモダリティ上に集合類似度  $sim(d_1, d_2)$  が分散とマッピング機能のために使用される。

10

## 【0131】

## マルチモードコレクション使用解析

ユーザをユーザの癖に従ってクラスタリング使用とすると困難がある。幾つかのケースにおいて、ウェブサイトのユーザをクラスタリングするために利用出来る唯一の直接情報は、ユーザが何れのページを如何なる頻度でアクセスするかということである。残念ながら、これは、類似度を決定するためには十分な情報がないので、しばしば相互に排他的なページビューでユーザをクラスタリングすることが不可能となる。

## 【0132】

このタイプの状況でマルチモードクラスタリングを可能とするために、媒介マルチモード表示がマトリックス操作によって計算される。例えば、 $P$  が  $n_p$  個のロー(ページの全数)と  $n_u$  個のコラム(ユーザの全数)を有するページアクセスのマトリックスであるとする。各コラムは、関数  $p$  によって発生されるベクトルに対応し、その微分は詳細に上述している。例えば、ユーザ番号 5 に対応する第 5 コラムは、 $p(u_5)$  である。 $T$  を  $n_p$  個のコラム(ページの全数)と  $n_t$  コラム(ワードの全数)を有するテキストマトリックスとする。上述のように、各コラムは関数  $t$  によって生成されるベクトルに対応する。例えば、ドキュメント番号 7 に対応する 7 番目のコラムは  $t(d_7)$  である。次に、ユーザのテキスト表示が以下のように計算される。

20

$$P_T = T \cdot P n_t$$

ローと  $n_u$  コラムを有するマトリックスであるこのマトリックスの内積は、各ユーザがアクセスしたページのテキストコンテンツの重み付け平均として解釈される。或いは、他の方法として、 $P_T$  は、アクセスされたページのコンテンツへのページアクセスの補外として解釈されてもよい。

30

## 【0133】

このアプローチの有用性の例として、個人の複写機 X C 5 4 0 を記述するページにアクセスした唯一のユーザの例を考える。モノモードクラスタリングがページアクセスに基づいてのみで実行される時、このユーザの類似度を他のユーザでアクセスすることは実際ではない。理由は、このユーザは、このページをアクセスした唯一のユーザのためである。また、ユーザは、積  $P_T = T \cdot P$  によって計算されるように、テキストモダリティに基づいて表される時、ユーザは、X C 5 4 0 ページ上に生じる"リーガルサイズ"又は"ペーパーレイ"のようなワードによって  $P_T$  に表される。このユーザのテキスト表示( $P_T$  においける単一のコラムによって定義されるベクトル)複写機ページをアクセスする他のユーザのテキスト表示に類似する。そして、上述のように、コサイン距離メトリックは、クラスタリング目的のために、 $P_T$  におけるユーザ間の類似度を決定するために使用され得る。この例は、媒介表示が類似度アセスメントとクラスタリングにおいて、どのように助けになるかを示す。

40

## 【0134】

更なる例によって、インリンク、アウトリンクと URL モダリティは、類似に計算される媒介によって表示可能である。ここで、マトリックス乗法は、 $L \cdot P$  (インリンク)、 $O \cdot P$  (アウトリンク)及び  $U \cdot P$  (URL) であり、ここで、 $L$ 、 $O$  及び  $U$  は、夫々

50

ンリンク、アウトリンク及びURLに対するマトリックスである。この概念は、ドキュメント当りベースに基づき計算されたあらゆる望ましいモダリティ又は特徴と共に、テキストジャンル、カラーヒストグラム及びカラー複雑さのような、その他のモダリティに対しても拡張され得る。

#### 【0135】

従って、ユーザがドキュメントコレクションとどのように対話するかを解析するためのマルチモード技術が可能である。この処理は、コレクションユーズ解析(CUA)と呼ばれる。ライブラリの編成及び解析において大きな文学があるが、これは、デジタルコレクションのための調査中の領域である。大部分の既知の従来の作業において、コレクションは、ユーザニーズの特徴付け(例えば、汎用クラスタリングによる)を行うこと無しに編成される。このセッションにおいて、実際のコレクション使用の解析は、コレクションの編成がどのように改良され得るか及びコレクションの何れの部分がユーザ母集団の特定のセグメントに対して最も価値があるかのような問題をどのように通報するかを示す。

10

#### 【0136】

これらの問題は、ウェブコレクションの豊富なハイパーリンク構造及びそれらの商業的重要性(これらの両方は、優れたコレクションデザインを必要とする)のために、ワールドワイドウェブのコンテキストにおいて特に重要である。表1(上記)にリストアップされたモダリティの内、以下の情報は、ページ及びユーザ:URL、アウトリンク、インリンク及び使用ログ、を特徴付けるために本発明の好適な実施の形態において使用される。上述のように、この情報の利用可能性は、CUAへのマルチモードアプローチに動機付けを行う。全ての可能なモダリティから利用出来る情報を利用し組合せることが出来ることが望ましい。

20

#### 【0137】

ここで記述されるCUAのために使用される主な技術は、ユーザのマルチモードクラスタリングであるが、これらのクラスタを解釈しようとする問題が残る。アブストラクトにおいて、クラスタのオブジェクトは、テキスト、使用、コレクショントポロジー(インリンクとアウトリンク)及びURLの特徴上のそれらのオブジェクト間の類似度によって特徴付けられる。オブジェクト間のこれらの特徴の類似度を明らかにするために、種々のユーザインターフェースと可視化技術が使用される。

30

#### 【0138】

ディスクツリー(後述される図23)は、ウェブサイトのページとハイパーリンクトポロジーを可視化するために使用されることが出来、且つ典型的にユーザの種々のクラスタに興味を持たせるサイトの部分を識別するのに有利であることが判った。また、ユーザのクラスタの興味の特徴を表すテキストとURLを要約するための技術は、本発明によって使用される。このような技術を組合せることによって、解析者には、自動的に識別されたクラスタの興味を特徴付けるテキスト、トポロジー及びURLの識別が提供され得る。

#### 【0139】

以下に述べられる例の実行において使用されるテストベッドは、1998年5月の17日と18日の24時間でのゼロックスウェブサイト(<http://www.xerox.com>)の完全なスナップショットより成った。約6400人のユーザに対する丸1日の使用情報が収集された。ユーザは、ブラウザクッキーに基づいて識別された。更に、全体のテキストとハイパーリンクトポロジーが抽出された。スナップショット時に、サイトは、6000以上のHTMLページと8000の非HTMLドキュメントから成った。

40

#### 【0140】

テストベッドは、3つの主要な要素から成った。即ち、モード情報を実数ベクトル( $R_n$ に埋め込まれた)にマッピングするマッピングプログラム、複数のセットのユーザをクラスタリングするクラスタリングプログラム及びウェブサイトの対話データ可視化を処理する可視化プログラムである。可視化プログラムは、ウェブサイトのディレクトリ構造を解析して、図23に示されるように、ディスクツリーを構成出来た。図示されているように、ウェブサイトの各ディレクトリは、全てのサブディレクトリを有するツリーの一つのノ

50

ードに対応し、ディレクトリ中のファイルは、そのノードの子として表される。好ましくは、そのツリーのレイアウトは、幅優先ベースで実行される。

【0141】

従って、本発明の実施の形態において使用される可視化システムは、図23に示されるように、ディスクツリーを構成し、ウェブサイトのベーシクトポロジーを表す。各ディレクトリは、一つのノードの子達として表されるディレクトリ中の全てのサブディレクトリとファイルを有するツリー中のそのノードに対応する。このツリーのレイアウトは、幅優先ベースで実行される。図23中のディスクツリー2310は、ゼロックス"スプラッシュページ" (<http://www.xerox.com/>) から開始する、ゼロックスウェブサイトを示し、引き続くディレクトリは、ディスクの中心から延出する同心リンクとして描かれる。これは、非対称ディスクを生成する。

10

【0142】

このディスクツリーは、アナリストユーザにクラスタについてのトポロジー情報へアクセスする方法を提供する。

【0143】

ディスクツリー2310において、クラスタは、そのクラスタの数に対応する全てのセグメントを一つのカラーに着色することによって可視化される。例えば、本発明の好適な実施の形態において、クラスタ中のメンバーシップは、そのクラスタ中のドキュメントに対応するセグメント2312、2314及び2316をレッド(図中では、太線で指示されている)で着色されることによって指示され得る。更に、好ましいシステムは、複数のメンバーシップの可視化を可能とする。これらの場合、複数のメンバーシップは、例えば、"レッドクラスタ"と"ブルークラスタ"中の同時メンバーシップを指示するためにレッドとブルーのストライプでセグメントの一つのグループ2320を着色することによって、そのページが属する全てのクラスタのカラーをミックスすることによって簡単に指示される。

20

【0144】

また、ダイアログボックスインターフェース(図24)を介して、本発明の好適な実施の形態のユーザは、何れのクラスタを表示するかを対話的に指定できる(本例では、同時に一つ又は二つのクラスタに制限される)。このダイアログボックスは、各クラスタのメンバーのテキスト表示をディスプレイする。各クラスタメンバーに対して、各モダリティの重みがリストアップされる。インリンク、アウトリンク、テキスト及び使用モダリティは、等しく重み付けされる(各々25%)。"クラスタリングレポート"2410は、ユーザクラスタに対する全てのドキュメントに亘って最も特徴的なキーワード2412を含む。これによって、他の特性を見ると同時にこのモダリティのハイレベルアブストラクションへの素早いアクセスが可能となる。"ドキュメントレポート"2414は、URLとそのクラスタ中の最も特徴的なドキュメントのテキストサマリー2418を提供する。多次元クラスタリングを有する経験は、幾つの場合、クラスタリングレポートがそのクラスタの最良の特徴付けであり、また他の場合、ドキュメントレポートが最良の特徴付けを提供する。そのシステムとの対は、全体のクラスタ又はその最も代表的なドキュメントのサマリー又はその両方を容易にアクセス出来ることによって多いに促進される。

30

40

【0145】

マルチモードクラスタリングの結果は、各モダリティに対するクラスタの最も特徴的な次元のテキストのリストアップである。例えば、クラスタが、Xerox Home Centre product (ゼロックスホームセンタープロダクト"について(about)")の場合、テキストモダリティに対する顕著な次元は、ワード"Home centre (ホームセンター)"である。それがテキストベッドゼロックスウェブサイトを与えられると、20から50個のクラスタが生成した。各クラスタは、数百のユーザを含み、テキストフォームのクラスタ結果の識別、比較及び評価のタスクは驚くべきものである。この場合、ディスクツリー(上述の)は有用である。

【0146】

50

図 2 4 に示されるように、クラスタレポートウィンドウ 2 4 1 0 は、ユーザクラスタに対する全てのドキュメントに及ぶ特徴的キーワード 2 4 1 2 を含む。これらは、クラスタのテキスト図心（テキストベクトルが図心を表す）中の最も高く重み付けされたワードを選択することによって計算される。このようなサマリーは、ユーザに大きなクラスタのテキストの信頼できるアセスメントを提供することが判った。

【 0 1 4 7 】

ドキュメントレポートウィンドウ 2 4 1 4 は、URL 2 4 1 6 と最も特徴的なドキュメント（クラスタのテキスト図心に最も近いドキュメント）のテキストサマリー 2 4 1 8 を提供する。それと共に、クラスタレポートウィンドウ 2 4 1 0 とドキュメントレポートウィンドウ 2 4 1 4 は、アナリストユーザに他のモダリティを見ると同時に、テキストモダリティのハイレベルアセスメントと URL を提供する。

10

【 0 1 4 8 】

図 2 4 のダイアログボックスインターフェースの残りの部分は、何れのクラスタが表示されるかを指定するために使用される。ダイアログボックスは、テキストを使用して、各クラスタのメンバーを表示する。各クラスタメンバーに対して、各モダリティの重み 2 4 2 0 がリストアップされ（本図に示されるクラスタリングは、五つのモダリティの内の四つに対して行われた）、本発明の好適な実施の形態においては、ユーザによって調節され得る。例えば、図 2 4 において、/ investor / pr / ir 9 8 0 5 1 2 . html は、クラスタ 0 のメンバーとして示される。インリンク、アウトリンク、テキスト及び使用モダリティは等しく重み付けされる（各 2 5 %）。

20

【 0 1 4 9 】

図 2 4 のダイアログボックスと図 2 3 のディスクツリーにユーザを直接示しの代わりにページを示す一つのモチベーションは、ユーザが、ページがされるのと同じ方法で、構造的且つ階層的に編成されない場合であり、その場合ユーザの直接可視化が困難となる。

【 0 1 5 0 】

従って、クラスタを提示するための二つの方法が提案される。第 1 の方法は、そのクラスタの全てのメンバーの可視提示から成る。上述のディスクツリーに基づき組み立てる場合、この方法は、メンバーが埋め込まれている階層構造がある場合、直進的である。例えば、ページのクラスタは、そのクラスタのメンバーに対応するディスクツリー中の全てのノードを着色することによって示される。

30

【 0 1 5 1 】

調停を介して表示されるオブジェクトのクラスタリングを等しく直進的に示す方法は無い。ディスクツリーとして可視化され得るユーザの直接階層的編成はない。従って、次に、技術的問題は、ウェブページベースの可視化においてユーザクラスタをどのように示すかである。この問題は、ランダムなユーザが望ましいクラスタから選択される場合、特定のページがアクセスされる確率を計算することによって解決される。確率  $P(p|u)$  は、ページ  $p$  がユーザ  $u$  によってアクセスされる相対頻度として計算される。例えば、ユーザが三つのページをアクセスする場合、その三つのページの各々は、 $1/3$  の確率  $P(p|u)$  を有する。次に、確率  $P(p|c)$ 、即ち、クラスタ  $c$  中の何れかのユーザによってページ  $p$  がアクセスされる相対頻度は、以下のように、そのクラスタ中のユーザに対する確率  $P(p|u)$  の平均として計算される。

40

【数 1 4】

$$P(p|c) = \sum_{u \in c} \frac{1}{|c|} P(p|u)$$

ここで、 $|c|$  はクラスタ  $c$  中のユーザの総数である。この可視化は、“密度プロット”と考えられ得る。直感的に、それは、このクラスタからの典型的なユーザが直面し得る問題に答える。本発明のこの好適な実施の形態において、全ての非 0 確率は、0.3 から 1.

50

0のスケールにマッピングされ、それによってクラスタ中のユーザによって2,3回しかアクセスされない偶数ページが明瞭に可視である。

【0152】

ユーザ母集団を解析するために、テストベッドの6400人のユーザ全てが20個のクラスタへクラスタリングされた。ユーザクラスタの内の9個は、ゼロックスプロダクト提供、例えば、P a g i s 走査、複写、X S o f t ソフトウエア、ゼロックスソフトウエアライブラリ(プログラムをダウンロードするため)、ホームとデスクトップ製品、及びウィンドウズ(登録商標)のためのT e x t B r i d g e、への興味によって特徴付けられた。7個のユーザクラスタは、一つの単一ページ、例えば、ドライバのインデックスやゼロックスホームページ、をアクセスした。ユーザの一つのクラスタは、使用情報をアクセスした。一つのクラスタは、ゼロックスについてのプレスリリースとニュースのような投資情報への興味によって特徴付けられた。二つのクラスタは、その他のカテゴリーに何れにもうまく当てはまらないユーザを含むようにミックスされた。従って、再び、図23を参照して、本発明の好適な実施の形態において、種々のセットのドキュメント2312、2314および2316は、カラーでハイライトされ、ユーザの特定のクラスタ(又は複数のクラスタ)がアクセスする見込みのあるドキュメントを指示する。

10

【0153】

クラスタを提示するための第2の方法において、テキストベースのクラスタサマリーは、各モダリティ毎に最も顕著な次元を提示することによって発生される。一例が、ゼロックスホームセンター(X e r o x H o m e C e n t r e)に興味のあるユーザのクラスタに対して表2に示されている。各モダリティ毎に、10個の最も顕著な次元がリストアップされる。10個の最も顕著なワード、このクラスタによってアクセスされたページを指す10個の最も顕著なページ、アクセスされたページ上に生じる10個の最も顕著なアウトリンク、アクセスされた10個の最も顕著なページ、及び10個の最も顕著なu r l エレメントである。クラスタ(この場合、ユーザ)中にあるオブジェクトのみに基づいてクラスタを解釈し比較することは、大変なタスクである。顕著な次元によるテキストサマリーは、クラスタ及びユーザが同じクラスタに置かれた理由を理解するのを一層容易にする。

20

【表 2】

テキスト		
0.504	8332	homecentre
0.221	14789	detachable
0.171	15270	artist
0.162	5372	slot
0.155	12010	mono
0.142	21335	photoenhancer
0.122	237	foot
0.121	4605	creative
0.113	3533	projects
0.109	21336	pictureworks

10

インリンク		
0.343	23856	products/dhc/index.htm
0.265	24144	products/dhc/06does.htm
0.259	17045	soho/whatsnew.html
0.257	24155	products/dhc/13inclu.htm
0.240	24151	products/dhc/07buser.htm
0.240	24152	products/dhc/07cuser.htm
0.235	24143	products/dhc/12more.htm
0.235	24157	products/dhc/15supp.htm
0.235	24156	products/dhc/14req.htm

アウトリンク		
0.527	24143	products/dhc/12more.htm
0.272	24156	products/dhc/14req.htm
0.272	24155	products/dhc/13inclu.htm
0.272	24157	products/dhc/15supp.htm
0.255	24149	products/dhc/11pagis.htm
0.248	31814	http://www.teamrx.com/retailers.html
0.216	24145	products/dhc/07user.htm
0.216	24144	products/dhc/06does.htm
0.192	23856	products/dhc/index.htm
0.137	23857	products/dwc450c/index.htm

20

ページ		
0.557	37067	products/dhc
0.330	24143	products/dhc/12more.htm
0.303	19452	products/multiprd.htm
0.287	24144	products/dhc/06does.htm
0.274	24739	soho/dhc.html
0.233	24155	products/dhc/13inclu.htm
0.208	24156	products/dhc/14req.htm
0.191	24148	products/dhc/09scan.htm
0.184	24157	products/dhc/15supp.htm
0.176	24145	products/dhc/07user.htm

30

url		
0.791	15	products
0.583	2036	dhc
0.141	646	soho
0.057	2037	dwc450c
0.054	895	print
0.044	31	cgi-bin
0.042	603	supplies
0.036	1768	usa
0.027	91	xps
0.020	844	wwwwais

40

## 【 0 1 5 4 】

所与のモダリティに対する顕著な次元は、 $P(p|c)$ によって表される確率を使用して集合特徴ベクトルへ寄与するドキュメントを重み付けることによって計算される。次に、集合特徴ベクトル中の最も大きな項が、顕著な次元を表す。例えば、表2を参照すると、図示されたクラスタに対する集合テキスト特徴ベクトルは、ワード"homecentre"に対応し、2番目に大きな項は、ワード"detachable"に対応する。集合

50

URL特徴ベクトルでは、最も重要なワードは、"products"であり、"dhc"が続く。

【0155】

アクセスされるコレクションの部分のこのような詳細な特徴付けは、適切なマテリアルの追加又は既存のマテリアルの改良のために使用され得る。例えば、唯一つの小さな投資家クラスタがあることを決定することは、驚くべきことである。これは、サイトに十分な投資情報が無い場合やそのレイアウトがより魅力的にするために改良されるべきである場合の証拠として解釈される。

【0156】

上述のように、幾つかのクラスタの顕著な特徴は、それらが唯一ページのみをアクセスするユーザからなることである。一例は、Text Bridge Pro 98 (光学的文字認識プログラム)のトライアルバージョンを要求するためのページのみにアクセスするクラスタである。これらのユーザは、明確に定義された情報の必要性があり、多分外部からのリンクに続いている。ユーザは、ユーザが必要とする情報(例えば、ゼロックスホームページ上のゼロックス社の株価)を得ると、直ちにユーザは離れる。

【0157】

他のクラスタは、俯角動作、ユーザが多数のページを介してブラウジングする時に、徐々に満足される多くのよりアモルファス情報の必要性によって特徴付けられる。一例は、小さなオフィスやホームオフィスに相応しいより小さなデバイス上で情報を得るドキュメントホームセンター(Document Homecentre)と呼ばれるサブ階層をブラウジングするユーザのクラスタである。経験的な解析において、このクラスタからのユーザが一般に幾つかの異なるドキュメントホームセンター製品に対応するサブ階層の幾つかのページを見るのが判った。明らかのように、これらのユーザは、ゼロックスウェブサイトに入り、利用出来る製品の範囲及び比較的広い範囲の情報を見ることを求める処理について学習する。

【0158】

コレクションの使用のこの解析は、再びより良いデザインへ反映され得る。例えば、しばしば一緒にブラウジングされる一セットのページは、ブラウジングを容易にするために、ハイパーリンクを介して一緒にリンクされるべきである。

【0159】

また、マルチモードユーザクラスタリングは、ウェブサイトのデザインを改良するために有用である。図23のディスクツリー2310は、50-クラスタのクラスタリングからの投資家のクラスタを示す。太い領域2312と2314によって指示される図面の上半分に強力なアクティビティの二つの領域がある。一つの領域2312は、サブ階層"annual report (年報)"に対応し、他の領域2314は、サブ階層"fact book (ファクトブック)"に対応する。多くの投資家が両方を見るという事実は、コレクションが再編成されたて、これら二つのサブ階層が共に検出されるべきであることを提案する。

【0160】

本システムは、探索データ解析のためにマルチモードクラスタリングを使用することの例である。本システムは、1998年5月17日のユーザ母集団を特徴付けるために使用された。6400人のユーザの全員が20個のクラスタに割り当てられた。9個のクラスタは、製品カテゴリーに対応する。Pags 走査、複写機、XSoftソフトウェア、ゼロックスソフトウェアライブラリ(ページをダウンロードするため)、ホームとデスクトップ製品、ウィンドウズのためのText Bridgeがある。7個のクラスタは、主に単一ページ、例えば、ドライバのインデックスやゼロックス社のホームページ、にアクセスするユーザに対応する。一つのクラスタは、使用情報をアクセスする投資家を含む。一つのクラスタは、投資家及びプレスリリースとゼロックス社についての他のニュースに興味がある他の投資家を含む。二つのクラスタは、その他のカテゴリーの何れにもうまく当てはまらないユーザを含む。このように、マルチモードクラスタリングは、アナリスト

10

20

30

40

50

がユーザ母集団の迅速な特徴付けを得る事を可能とする。

【0161】

ディスクツリーを含む多くの可視化は、スクリーン上に制限された数のノードのみを描くことが出来る。マルチモードクラスタリングは、ノードのグループをメタノードへのノード集合のために使用され得る。例えば、スクリーンのエッジ上にノードの1000個のサブノードをディスプレイするのに必要なスクリーン領域が無い場合、これら1000個のサブノードは、マルチモードクラスタリングを使用して、5個のメタノードへ集合され得る。次に、5個のメタノードのディスプレイは、全1000個のサブノードをディスプレイするよりも小さい空間を取る。

【0162】

また、マルチモードクラスタリングは、データ最小化のために使用され得る。ユーザのクラスタがマルチモードアルゴリズムによって生成されると、ユーザは、顕著な特徴を自動的に検出することが出来る。例えば、顕著なワードとして"homecentre"を示す表2のHomeCentreクラスタのテキスト表示に基づいて、ユーザは、如何に上手に"homecentre"単独によって特徴付けられるかをテスト出来る。

【0163】

他の一つのデータ最小化アプリケーションは、例外的なオブジェクトの発見である。例えば、訴訟の開示段階において、法律事務所は、型通りの記事を大部分含む類似のドキュメントの大きなグループにではなくて、孤立するドキュメントに興味を示すに過ぎないかもしれない。マルチモードクラスタリングは、大きなグループの類似ドキュメント(例えば、共有する型通りの記事)を識別する。次に、興味あるドキュメントは、大きなクラスタの図心から大きく離れたもの間にあるかもしれない。

【0164】

本発明のデータ最小化技術は、第1のグループに対してマルチモードクラスタリングを実行し、次に第2のグループを第1のグループのクラスタへ割り当てることによって二つのグループのオブジェクトを比較する。この解析技術は、ウェブサイトのゼロックススペースのユーザと非ゼロックススペースのユーザを比較するのに成功裏に使用され、ゼロックス従業員がゼロックス製品のユーザでありためにほんの僅かな差が発見された。それは、外部ゼロックスウェブサイト(ドライバをダウンロードするため、製品情報を参照するため等)へ行く主な理由の一つである。一つの差は、より高い比率のゼロックスユーザがページのみ、ゼロックスホームページ、を訪れることであった。その理由は、多分ゼロックス従業員の多くのブラウザが彼等/彼女等のデフォルトページとしてゼロックスホームページを有することであるので、ユーザは、ユーザのブラウザを開始すると自動的にゼロックスホームページに行き、次に異なるサイトのゼロックスホームページに行く。この例は、異なるユーザグループを比較するために、マルチモードクラスタリングのユーティリティを証明する。

【0165】

イントラネットを含む、大きなコレクションを編成するための益々重要な技術は、階層的クラスタリングである。その目的は、ヤフー上で(及び多くのイントラネット上で)見つけられ得るように、階層を自動的に発生することである。階層的マルチモードクラスタリングは、そのような階層を自動的に発生するために或いは類別する人に手作業の編集を可能する第1のカットを与えるために使用され得る。

【0166】

コレクション使用解析の基づく推奨  
最後に、マルチモードユーザクラスタに基づく推奨システムは、上述されたように、マルチモードコレクション使用データのコレクションを有することが可能である。一セットのクラスタは、一トレーニングセットのユーザから導かれる。新たなユーザは、2, 3回の初期ページアクセスに基づいてクラスタの一つへ割り当てられる。次に、割り当てられたクラスタ中のユーザによってアクセスされたページは、ユーザに推奨される。クラスタリングがマルチモード情報に基づいてなされるので、それは、有用な推奨を行うのに十分に

10

20

30

40

50

頑強である。

【0167】

本発明のマルチモード推奨システムは、図25に示されている。最初に、ユーザのトレーニングセットが識別される(ステップ2510)。ユーザについて利用可能なあらゆるタイプの情報が収集される。ここで開示された実施の形態において、ページユーザアクセスに関する情報、並びにテキストコンテンツ、インリンク、アウトリンク、及びこれらのページのURLを収集するために有用であることが判った。また、リアルタイムドキュメントアクセスがこの(そのデータは使用ログから又は利用出来る時ユーザのセットのブラウザ"bookmarks"ですらからもであってもよい)ために使用される必要は無いことに注意すべきである。また、上述から気付かれるように、デモグラフィック情報及び他の種類の追跡された情報のような、このアプリケーションにおいて有用であり得るユーザへ適用可能な他のモダリティ(ページ使用を超えて)がある。

10

【0168】

次に、ユーザは、マルチモードクラスタリングに関連するセクションにおいて、上述のように、マルチモード情報を介してクラスタリングされる(ステップ2512)。ページ使用が、本発明の好適な実施の形態におけるように、ユーザについて収集された主情報である場合、上述のように、種々のドキュメント特徴によってユーザの調停された表示を介してユーザをクラスタリングすることが適切である。他の戦略のまた可能であることを認識すべきである。例えば、デモグラフィック情報が収集されると、デモグラフィック情報に関して簡単にユーザをクラスタリングすることがより適切である。クラスタリングされるベースの選択は、本発明のシステムのデザイナーの判断に任される実施の詳細事項である。或いは、その選択は、ユーザに任されてもよい。

20

【0169】

新たなユーザが無い場合(ステップ3514)、その処理は終了される(ステップ2516)。或いは、新たなユーザが識別されると(ステップ2518)、ブラウジング情報が新たなユーザから収集され(ステップ2520)、ユーザが最も近い既存のクラスタに割り当てられる(ステップ2522)。本発明の好適な実施の形態において、ユーザは、上述のように、テキストコンテンツ、インリンク、アウトリンク及びURLに亘って計算された集合コサイン類似度に基づいて、割り当てられる。

【0170】

次に、最も近いクラスタ中の最もポピュラーなページが識別され(ステップ2514)、新たなユーザへ推奨される(ステップ2526)。本発明の他の実施の形態において、ネーム、メールアドレス、又は最も近いクラスタにおけるユーザのための他の識別情報(或いは、上述の集合コサイン類似度メトリックを介して識別される、その最も近いクラスタ中の少なくとも一つのユーザ)が、その新たなユーザに提供され、それによってその新たなユーザが望ましい領域中に"エキスパート"を識別出来る。

30

【0171】

このアルゴリズムは、他の推奨アルゴリズム以上に幾つかの利点を有する。このアルゴリズムは高速である。クラスタリングがコンパイルタイム動作であるので、唯一のランタイム動作がマルチモード情報の各モダリティのベクトル空間へのマッピングと各クラスタを有する集合コサイン類似度の演算である。これは効率的である。同じ利点を得る他の方法は、ユーザ母集団をサマリー化する一つの方法としてクラスタリングを見なすことである。これは、ユーザ母集団が大きい場合に重要である。例えば、百万人のユーザの追跡を維持しなければならない代わりに、推奨は、唯、即ち1000人のユーザのみに基づいて、行われ、1000個のクラスタを代表するものは、完全なユーザ母集団から導出される。

40

【0172】

ユーザ母集団からクラスタを導出することは、新たなユーザを割り当てるよりのより高価であるが、1日に数回或いはそれを越える頻度で行われることはなお十分に効率的である(クラスタリングは、クラスタリングされるべきオブジェクトの数に関してリニアであ

50

るので)。このように、推奨は、ユーザのニーズを素早く変化するために採用し得る。これは、図26に示されるように実行される。そうすることが望ましい場合（例えば、定期的に或いは十分な数の新たなユーザがユーザプールに追加された場合）、サブセットのユーザが初めに識別される（ステップ2610）。上述のように、大きな母集団では、サブセットのユーザは、全体の母集団の特徴を非常に良好に表示出来る。次に、そのセットのユーザは、再クラスタリングされる（ステップ2612）。次に、各ユーザ毎に最もポピュラーなページが決定され（ステップ2614）、従って新たなユーザに推奨されたページが調節される（ステップ2616）。

【0173】

コレクション使用解析に基づくマルチモード推奨のためのここで記述されるアルゴリズムは、非常に正確で頑強であることが判った。他の推奨アルゴリズムは、新たなユーザの前のユーザとの比較に依存する。推奨がたまたま最も近い隣人である一人又は二人のユーザに基づく場合、孤立するものが推奨されたページに影響を及ぼし得るので、間違っただページが推奨されるかもしれない。クラスタベースの汎用化は、孤立するものの影響を減少する。更に、全ての利用可能情報が使用され且つ組み合わせられるので、そのアルゴリズムは、情報の単一源に依存する推奨アルゴリズムよりも頑強である。

【0174】

以下に記述される例に対して、テストベッドユーザ（即ち、1998年5月17日及び18日のゼロックスウェブサイトをユーザ）のアクションがログされた。ユーザのブラウジング癖に基づいて、これらのユーザが200個のクラスタに配された。

【0175】

クラスタベースのシステムによって行われる推奨の第1のタイプが表3に示されている。

【表3】

クラスタ35		
0.976277 probsum		
16406	0.088639	products/copiers.htm
37005	0.085385	http://www.xerox.com
19453	0.059099	products/cop_soho.htm
33739	0.051071	soho/xc0355.html
21231	0.040836	soho/xc1044.html
17033	0.039741	soho/xc0830.html
37025	0.036496	cgi-bin/wwwwais
19451	0.035938	products/cop_pers.htm
17029	0.034706	soho/xc0540.html
17010	0.028586	soho/5306.html
21232	0.026014	soho/xc1045.html

表3は、確率  $P(p|u)$ （上記を参照：ページ  $p$  の確率は我々がクラスタ35からのユーザ  $u$  を有することが与えられる）の計算に基づいて、ユーザクラスタ35に対する最もポピュラーなページを示す。この情報は、ページ "products/copiers.htm" をアクセスするあらゆるユーザへそのクラスタ中のその他のページ、即ち最もポピュラーな複写機、を推奨することによって宣伝され得る。これらのリンクの幾つかは、ページ "products/copiers.htm" からアクセス可能である。このアルゴリズムは、ユーザが最も関連する可能性のあるリンクを選択することを容易とする。

【0176】

クラスタベースの汎用化によって可能とされる第2のタイプの推奨が表4に示されてい

る。

【表 4】

クラスタ127		
1.000000 probsum		
24663	0.297222	employment/ressend.htm
37057	0.268162	employment
24666	0.079701	employment/resascii.htm
21384	0.076923	research/xrcc/jobopps.htm
37005	0.054701	http://www.xerox.com
37087	0.050000	cgi-bin/employment/xrxresume.cgi
24675	0.047436	employment/restip.htm
24664	0.023077	employment/college.htm
15355	0.012821	XBS/employt.htm
24665	0.012821	employment/recruit.htm
34418	0.012821	employment/overview.htm
37025	0.012821	cgi-bin/wwwais

10

20

この表は、ユーザクラスタ127に対する最も顕著なページを含む。このクラスタのコンテンツに基づいて、本システムは、種々の細分割のemployment（雇用）ページをジョブのために容易に適用するユーザへ推奨することが出来る。リストアップされたドキュメントは、中央のemploymentページ（数値識別子37057を有するその表の第2ページ）から直接にはアクセス出来ないゼロックスのウェブサイト上の幾つかのemploymentページを含む。二つのこのような直接にはアクセス出来ないページは"research/xrcc/jobopps.htm"と"XBS/employt.htm"である。このタイプの推奨によって、ユーザはユーザがそうでない場合（丁度時間を節約することとは反対である）は全く発見できないかもしれない何かを発見することが出来る。上述されたものと同じアルゴリズムは、これを、即ち、（幾つかの初期ページアクセス後の）新たなユーザをユーザに割り当てそしてそのユーザがアクセスしなかったクラスタのページ特性を推奨すること、完成するために使用される。

30

【0177】

表5は、ユーザクラスタ25に対する最も顕著なページを含む。

【表 5】

クラスタ25		
0.998387 probsum		
37057	0.661425	employment
37005	0.300403	http://www.xerox.com
34418	0.022581	employment/overview.htm
12839	0.004435	searchform.html
24675	0.004032	employment/restip.htm
37155	0.002688	scansoft/tbapi
37113	0.002016	factbook/1997
23465	0.000806	xbs

40

これらのユーザは、ブラウジングしており、多分ジョブへの適用は容易ではなく、XBS

50

のような指定の分割の employment ページはユーザに推奨されない。表 4 と表 5 との間のコントラストは、マルチモードクラスタリングによって発見される汎用化の例である。第 1 のクラスタのユーザは、ユーザの概要を提出することが多いにあり得る。XBS のような細分割の employment (雇用) ページをユーザに推奨することは優れたアイデアである。理由は、ユーザは、ジョブを見つけることについて厳しいと思うからである。

【0178】

他方、第 2 のクラスタのユーザは、何らかの一般的ブラウジングを行う。employment は、ユーザのブラウジングの目的であり、ユーザは、目的のジョブサーチを実行するようには思わない。これらのユーザは、ジョブの広告を有するページを見ようとはあまりせず、従って細分割の employment ページが彼等/彼女等には推奨されない。

10

【0179】

本明細書中に記載されている 3 つの符号、aaa, bbb, ccc は、下記の通り、便宜上置きかえたものである。

【外 1】

【0180】

aaa →  $\bar{c}$

bbb →  $\bar{x}_i$

ccc →  $\bar{x}'_i$

20

【図面の簡単な説明】

【0181】

【図 1】本発明に従うシステムとの使用に適するネットワークに接続されたドキュメントコレクションを示すブロック図である。

【図 2】コレクションに追加された新たなドキュメントを処理するために本発明の一実施の形態によって使用される処理を示すフロー図である。

【図 3】種々のドキュメントとユーザを表す特徴ベクトルを計算するために本発明の一実施の形態によって使用される処理を示すフロー図である。

30

【図 4】本発明の一実施の形態において、テキストベースの特徴ベクトルを計算するために使用される処理を示すフロー図である。

【図 5】本発明の一実施の形態において、テキストジャンル特徴ベクトルを計算するために使用される処理を示すフロー図である。

【図 6】本発明の一実施の形態において、カラーヒストグラム特徴ベクトルを計算するために使用される処理を示すフロー図である。

【図 7】本発明の一実施の形態において、対応する対のカラー複雑さ特徴ベクトルを計算するために使用される処理を示すフロー図である。

【図 8】本発明の一実施の形態において、ページ使用ベクトルを計算するために使用される処理を示すフロー図である。

40

【図 9】本発明の一実施の形態において、初期クラスタセンタを識別するためにウェブフロントクラスタリングで使用される処理を示すフロー図である。

【図 10】関連するオブジェクトをクラスタに割り当てるために k 平均クラスタリングで使用される処理を示すフロー図である。

【図 11】異なるモダリティにおけるコレクションオブジェクトの分散収集の仮説セッションを示す図である。

【図 12】照会 "ancient cathedral (昔の大聖堂)" に応答して戻されるテキストクラスタの例示の可視ディスプレイである。

【図 13】図 12 の第 1 のテキストクラスタを分散した後に戻されるテキストクラスタの

50

例示の可視ディスプレイである。

【図14】複雑さ特徴に基づくクラスタリングの後に戻される画像クラスタの例示の可視ディスプレイである。

【図15】照会"paper money (紙幣)"に応答して戻されるテキストクラスタの例示の可視ディスプレイである。

【図16】複雑さ特徴に基づく図15の第1のテキストクラスタのクラスタリングの後に戻される画像クラスタの例示の可視ディスプレイである。

【図17】カラーヒストグラム特徴に基づく図16の第3及び第4画像クラスタをクラスタリングした後に戻される画像クラスタの例示の可視ディスプレイである。

【図18】カラーヒストグラム特徴に基づく図17の第2の画像クラスタをクラスタリングした後に戻される画像クラスタの例示の可視ディスプレイである。

【図19】照会"pyramid egypt (ピラミッドエジプト)"に応答して戻されたテキストクラスタの例示の可視ディスプレイである。

【図20】複雑さ特徴に基づくクラスタリングの後に戻される画像クラスタの例示の可視ディスプレイである。

【図21】カラーヒストグラム特徴に基づくクラスタリングの後に戻される画像クラスタの例示の可視ディスプレイである。

【図22】図21の画像のセットを拡張し、カラーヒストグラムの基づく結果をクラスタリングした後に戻されるテキストクラスタの例示の可視ディスプレイである。

【図23】本発明によるクラスタの例示の間接可視化であり、一つのユーザクラスタは、そのクラスタのメンバーによって選択される可能性の高い全てのページを赤で(ここでは、矢印によって指示される)着色することによって図示される。

【図24】本発明の一実施の形態において、クラスタとドキュメントのコンテンツをブラウジングし示すために使用されるインターフェースを示す例示の可視ディスプレイである。

【図25】本発明に従う例示の推奨システムにおいてポピュラーなページをユーザに推奨するために使用される処理を示すフロー図である。

【図26】本発明に従う例示の推奨システムにおいて推奨を再計算するために使用される処理を示すフロー図である。

【符号の説明】

【0182】

- 120 コレクション
- 110 ドキュメント
- 114 テキストベクトル
- 116 URLベクトル
- 112 特徴ベクトル
- 124 通信ネットワーク
- 122 プロセッサ
- 124 ネットワーク
- 126 データベース
- 128、130、132 ユーザ端末

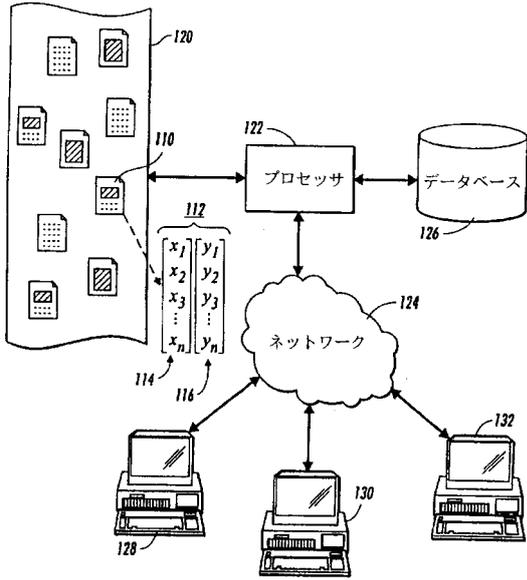
10

20

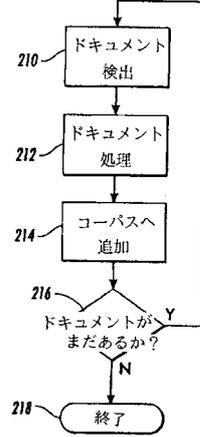
30

40

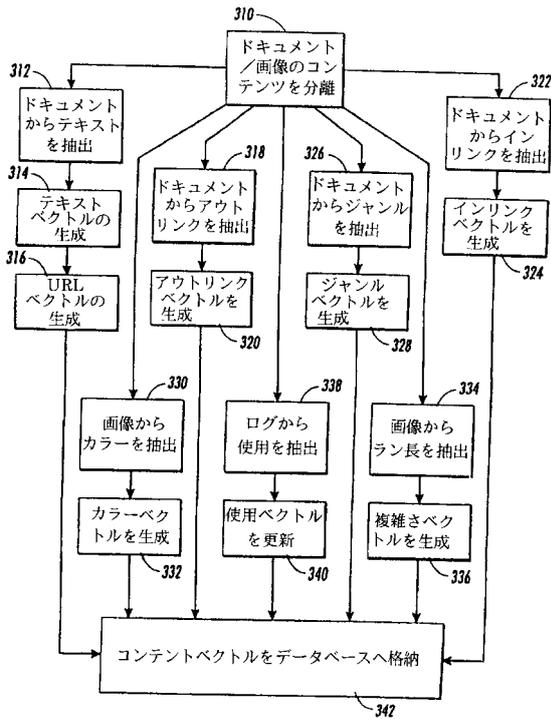
【図1】



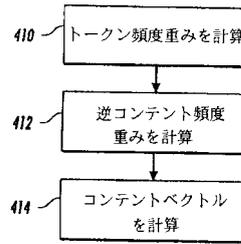
【図2】



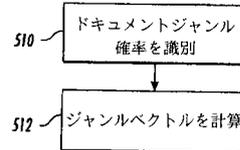
【図3】



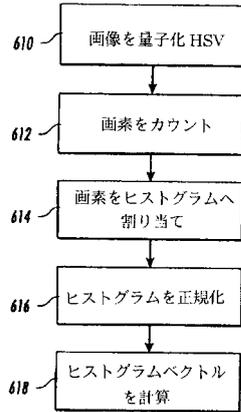
【図4】



【図5】



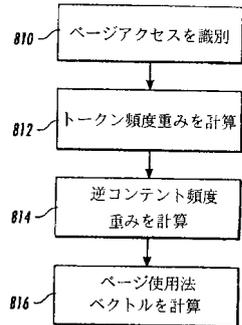
【図 6】



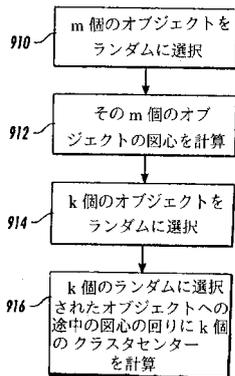
【図 7】



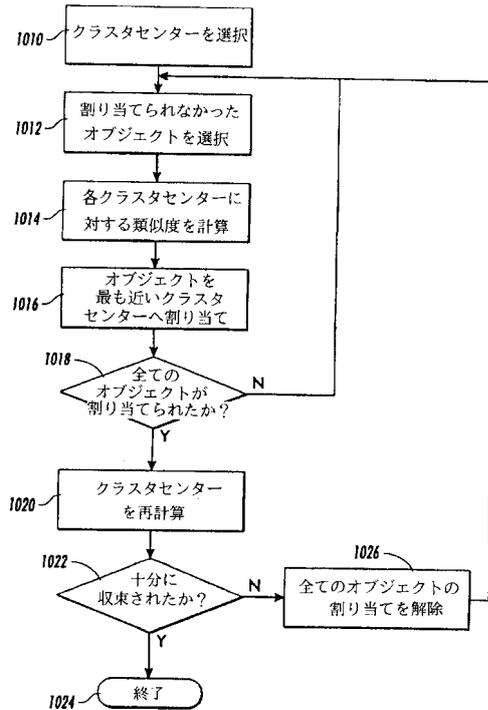
【図 8】



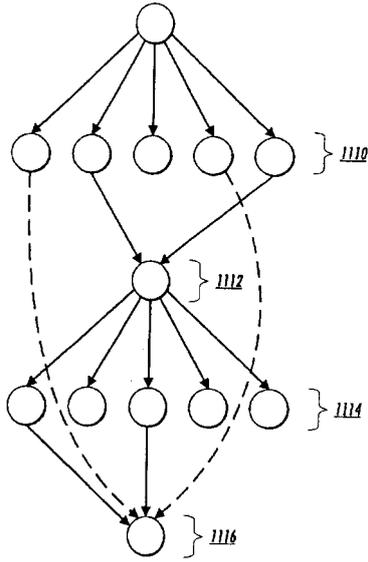
【図 9】



【図 10】



【図 1 1】



【図 1 2】

Figure 12 is a screenshot of a web browser interface. At the top, there is a search bar containing the text "ancient cathedral". Below the search bar, there is a list of search results, each with a checkbox and a small image icon. The results are labeled with numbers 1 through 5. To the right of the search results, there is a navigation panel with several buttons: "実行", "再入力", "次へ", and "URL". The interface also shows a search bar with the text "ancient cathedral" and a list of search results with checkboxes and small image icons.

【図 1 3】

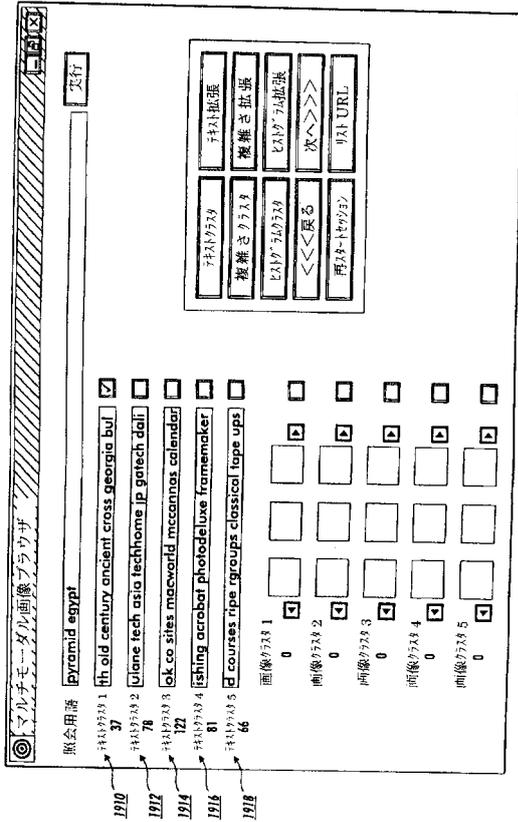
Figure 13 is a screenshot of a web browser interface. At the top, there is a search bar containing the text "ancient cathedral". Below the search bar, there is a list of search results, each with a checkbox and a small image icon. The results are labeled with numbers 1 through 5. To the right of the search results, there is a navigation panel with several buttons: "実行", "再入力", "次へ", and "URL". The interface also shows a search bar with the text "ancient cathedral" and a list of search results with checkboxes and small image icons.

【図 1 4】

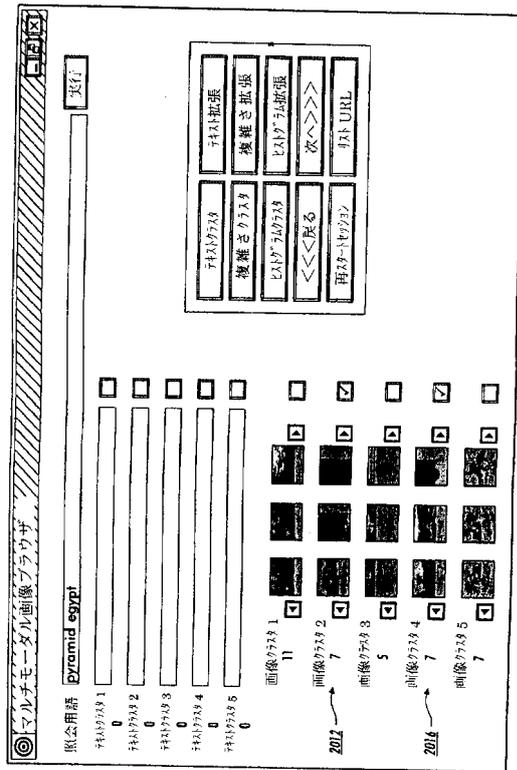
Figure 14 is a screenshot of a web browser interface. At the top, there is a search bar containing the text "ancient cathedral". Below the search bar, there is a list of search results, each with a checkbox and a small image icon. The results are labeled with numbers 1 through 5. To the right of the search results, there is a navigation panel with several buttons: "実行", "再入力", "次へ", and "URL". The interface also shows a search bar with the text "ancient cathedral" and a list of search results with checkboxes and small image icons.



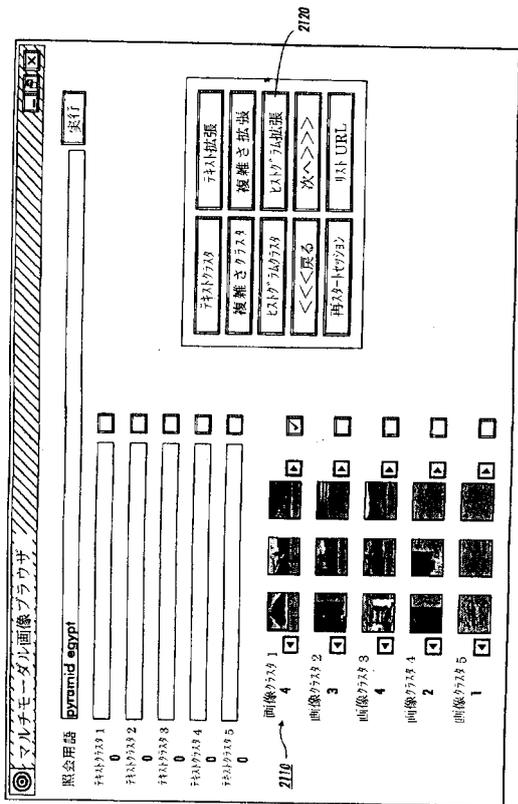
【 19 】



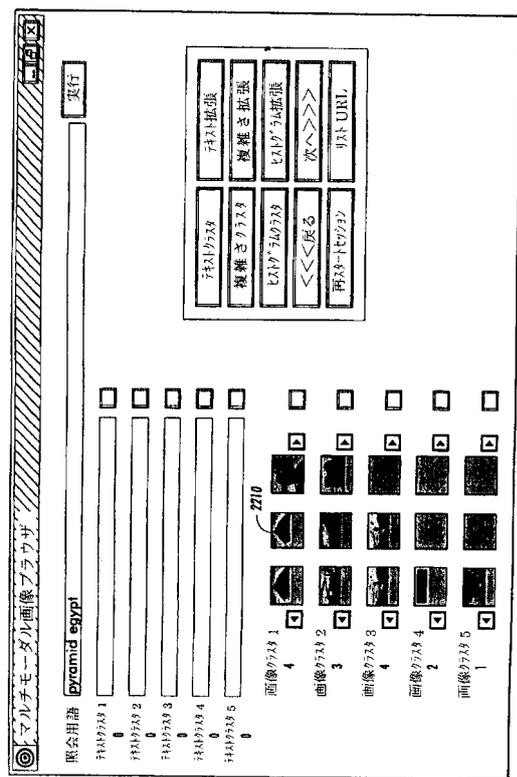
【 20 】



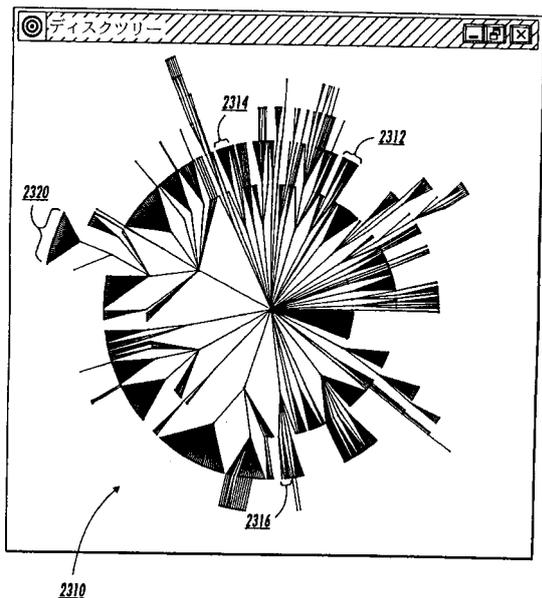
【 21 】



【 22 】



【図 23】



【図 24】

クラスターリング結果ビュー

クラスターリングレポート

クラスター ID: 0

トキメント (特設トキメントで表示): investor/results/96/196q/c

トキメントを指定: investor/pr/980512.htm ID: 36828

指定トキメントのコンテンツサマリー: xerox best products web site information service use documents printers copiers new today printer speed per pages document dual ppm

インスタラクション:

クラスタードキュメントのコンテンツサマリーを見る  
 ためにクラスタードキュメントはクラスタリング結果ビュー  
 中のトキメントメニューをダブルクリック。  
 又は編集ホックスヘムをキーピング

グラフ中の全てのクラスターを表示

グラフ中に示されるクラスターを指定

第1クラスター: [ ] 第2クラスター: [ ]

グラフ中にサブセットトリーを表示

グラフ中にリンクセンタービューを表示

グラフ上にラベルを表示

トキメントを見るためにブラウザを印刷

トキメントレポート

トキメントを指定: investor/pr/980512.htm

指定トキメントのコンテンツサマリー:  
 file items financial earnings per share growth revenues stock  
 consolidated operations common discontinued shares insurance  
 payments per director million executive annual loss shareholders  
 cash interest

サブセット中の  
 全トキメント数: 7581  
 全クラスター数: 5

モダリティ: 重み

リンク: 0.250000

リンク: 0.250000

セント: 0.250000

母田: 0.250000

2420

Cluster0/investor/pr/980512.htm
investor/pr/980512.htm
facebook/1998/contacts.htm
facebook/1998/orgchat.htm
facebook/1998/finance.htm
facebook/1998/historic.htm
facebook/1998/officers.htm
facebook/1998/other.htm
facebook/1998/balance.htm
facebook/1998/newenter.htm
facebook/1998/stock.htm
facebook/1998/groups.htm

クラスターリングレポート

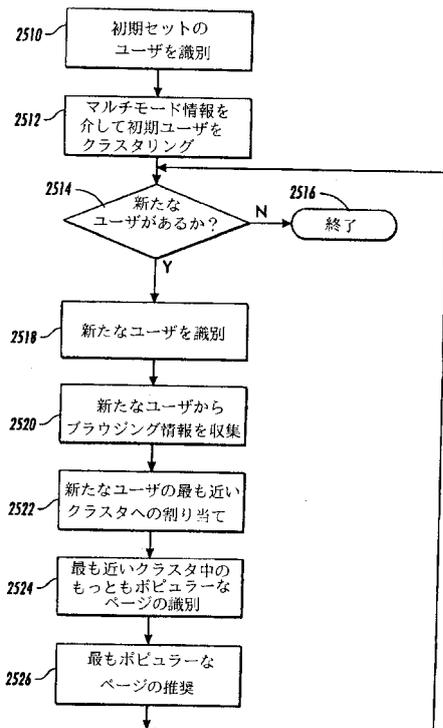
特徴

クラスター: 0 トキメント

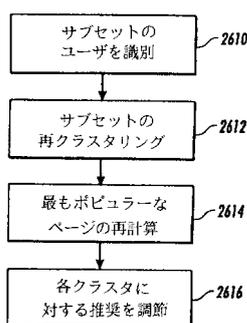
クラスター ID: investor/results/96/196q/c

クラスターのコンテンツサマリー  
 (特設トキメントで表示):  
 file items financial earnings per share growth revenues stock  
 consolidated operations common discontinued shares insurance  
 payments per director million executive annual loss shareholders  
 cash interest

【図 25】



【図 26】



## フロントページの続き

- (31)優先権主張番号 421419  
 (32)優先日 平成11年10月19日(1999.10.19)  
 (33)優先権主張国 米国(US)  
 (31)優先権主張番号 421767  
 (32)優先日 平成11年10月19日(1999.10.19)  
 (33)優先権主張国 米国(US)  
 (31)優先権主張番号 421770  
 (32)優先日 平成11年10月19日(1999.10.19)  
 (33)優先権主張国 米国(US)  
 (31)優先権主張番号 425038  
 (32)優先日 平成11年10月19日(1999.10.19)  
 (33)優先権主張国 米国(US)  
 (31)優先権主張番号 425039  
 (32)優先日 平成11年10月19日(1999.10.19)  
 (33)優先権主張国 米国(US)

- (72)発明者 ヒンリッヒ シュエツェ  
 アメリカ合衆国 9 4 1 3 1 - 1 5 5 2 カリフォルニア州 サンフランシスコ ポートラ ドライブ 1 0 0 ナンバー 1  
 (72)発明者 ウラス ガーギ  
 アメリカ合衆国 1 6 8 0 3 ペンシルベニア州 ステイト カレッジ ウェスト クリントン アベニュー 2 3 4  
 (72)発明者 ジェームズ イー・ピットカウ  
 アメリカ合衆国 9 4 3 0 6 カリフォルニア州 パロ アルト エルスワース プレイス 7 4 2  
 (72)発明者 ピーター エル・ピロリ  
 アメリカ合衆国 9 4 1 1 6 カリフォルニア州 サンフランシスコ スロート ブルバード 2 9 5 8  
 (72)発明者 エド エイチ・チ  
 アメリカ合衆国 5 5 4 1 7 - 1 9 3 7 ミネソタ州 ミネアポリス ショアビュー アベニュー サウス 5 2 4 1  
 (72)発明者 ジュン リ  
 アメリカ合衆国 4 7 4 0 1 インディアナ州 ブルーミントン イースト セカンド ストリート 2 1 0 6 ナンバー 7

審査官 鈴木 和樹

- (56)参考文献 特開平11-224256(JP,A)  
 特開平08-190578(JP,A)  
 米国特許第05835087(US,A)  
 吉田勝彦、外1名、インターネット上のマルチメディアデータの柔軟な検索方式の提案、電子情報通信学会技術研究報告(KBSE97-20~24)、日本、社団法人電子情報通信学会、1997年11月28日、第97巻、第413号、p.17-22  
 Rohini K. Srihari, Automatic Indexing and Content-Based Retrieval of Captioned images, Computer [online], 1995年9月、第28巻、第9号、p.49-56, [DL from IEEE Xplore]  
 Marco La Cascia、外2名、Combining Textual and Visual Cues for Content-based Image Retrieval on the World Wide Web, Content-Based Access of Image and Video Libraries, 1998.

Proceedings. IEEE Workshop on [online], 1998年 6月21日, p. 24 - 28, [DL from IEEE Xplore]

Ron Weiss、外6名, HyPursuit: A Hierarchical Network Search Engine that Exploits Content-Link Hypertext Clustering, HYPERTEXT '96 Proceedings of the seventh ACM conference on Hypertext [online], 米国, ACM, 1996年, p. 180 - 193, [DL from ACM Digital Library]

(58)調査した分野(Int.Cl., DB名)

G06F 17/30

IEEE Xplore

THE ACM DIGITAL LIBRARY