

(19) 日本国特許庁 (JP)

(12) 特 許 公 報 (B2)

(11) 特許番号

特許第4838529号
(P4838529)

(45) 発行日 平成23年12月14日 (2011.12.14)

(24) 登録日 平成23年10月7日 (2011.10.7)

(51) Int.Cl. F I
G 0 6 F 17/30 (2006.01)
 G 0 6 F 17/30 2 1 0 D
 G 0 6 F 17/30 1 7 0 A

請求項の数 22 外国語出願 (全 30 頁)

(21) 出願番号	特願2005-118052 (P2005-118052)	(73) 特許権者	500046438
(22) 出願日	平成17年4月15日 (2005.4.15)		マイクロソフト コーポレーション
(65) 公開番号	特開2005-302043 (P2005-302043A)		アメリカ合衆国 ワシントン州 9805
(43) 公開日	平成17年10月27日 (2005.10.27)		2-6399 レッドモンド ワン マイ
審査請求日	平成20年4月15日 (2008.4.15)		クロソフト ウェイ
(31) 優先権主張番号	10/826, 159	(74) 代理人	100077481
(32) 優先日	平成16年4月15日 (2004.4.15)		弁理士 谷 義一
(33) 優先権主張国	米国 (US)	(74) 代理人	100088915
			弁理士 阿部 和夫
		(72) 発明者	ベンユー チャン
			アメリカ合衆国 98052 ワシントン
			州 レッドモンド ワン マイクロソフト
			ウェイ マイクロソフト コーポレーシ
			ョン内

最終頁に続く

(54) 【発明の名称】 検索語提案のためのマルチ型データオブジェクトの強化されたクラスタリング

(57) 【特許請求の範囲】

【請求項 1】

コンピューティングデバイスにより実行される方法であって、

第1の型の少なくとも1つのオブジェクト、および前記第1の型とは異なる第2の型の少なくとも1つのオブジェクトを含むマルチ型データオブジェクトの間の関係を処理ユニットが決定するステップであって、前記関係は、前記マルチ型データオブジェクトの間のレイヤ内関係か、レイヤ間関係かのうち少なくとも1である、ステップと、

前記関係に基づいて前記マルチ型データオブジェクトを繰り返しクラスタリングして強化されたクラスタを前記処理ユニットが生成するステップと、

ユーザから送信される語に関係のある、提案される検索語のリストを、前記強化されたクラスタを利用して前記処理ユニットが生成するステップであって、前記検索語は、前記ユーザから送信される語の受信に回答して生成される、ステップと、

ネットワークインターフェースを介して、前記提案される検索語のリストをユーザへ前記処理ユニットが送信するステップと、

次の

【数 1】

$$\left\{ \begin{array}{l} a(X) = \beta L_X^T h(X) + (1 - \beta) L_{XY} i(Y) \\ h(X) = \beta L_X a(X) + (1 - \beta) L_{XY} i(Y) \\ i(X) = a(X) + h(X) \\ \\ a(Y) = \gamma L_Y^T h(Y) + (1 - \gamma) L_{YX} i(X) \\ h(Y) = \gamma L_Y a(Y) + (1 - \gamma) L_{YX} i(X) \\ i(Y) = a(Y) + h(Y) \end{array} \right.$$

10

を使用してノードの権威スコアおよびハブスコアを更新することによって、オブジェクト型内および異なるオブジェクト型の間での前記マルチ型データオブジェクトの個々のオブジェクト重要度の相互強化を前記処理ユニットが行うステップであって、 $X = \{x_1, x_2, \dots, x_m\}$ および $Y = \{y_1, y_2, \dots, y_n\}$ は R_X 、 R_Y 、 R_{XY} 、および R_{YX} の関係を有する異質のオブジェクト型のそれぞれのオブジェクトの集合を表現し、方向性が考慮される場合、 L_X および L_Y はそれぞれ集合 X および Y 内の関係を識別するリンクの隣接する行列を表現し、 L_{XY} および L_{YX} は X 中のオブジェクトから Y 中のオブジェクトへの関係を識別するリンクの隣接する行列を表現し、 $a(X)$ および $h(X)$ はそれぞれ X 内のノードの権威スコアおよびハブスコアであり、 $a(Y)$ および $h(Y)$ は、 Y 内のノードの権威スコアおよびハブスコアを表し、 $i(X)$ および $i(Y)$ はそれぞれ X および Y 中のノードの重要度を表し、およびは異なる関係に由来するリンクの影響を調整するための重みパラメータである、ステップと

20

を備えたことを特徴とする方法。

【請求項 2】

前記レイヤ間関係は、コンテンツに関係のある情報、関連トピックに対するユーザの関心、および関連 Web ページに対するユーザの関心のうち少なくとも 1 つを含むことを特徴とする請求項 1 に記載の方法。

【請求項 3】

前記レイヤ内関係は、クエリ改良、推奨される Web ページ、およびそれぞれのユーザの間の関係のうち少なくとも 1 つを含むことを特徴とする請求項 1 に記載の方法。

30

【請求項 4】

前記マルチ型データオブジェクトの各々は、選択された Web ページ型およびユーザ情報型のうち少なくとも 1 つに関係のあることを特徴とする請求項 1 に記載の方法。

【請求項 5】

前記レイヤ内関係は、前記マルチ型データオブジェクトの関連付けるオブジェクトに対する重要度を示すための第 1 の重み付けの仕組み、および前記第 1 の重み付けの仕組みとは異なる第 2 の重み付けの仕組みを含むことを特徴とする請求項 1 に記載の方法。

【請求項 6】

前記識別するステップおよび繰り返しクラスタリングするステップは、検索語提案のために実行されることを特徴とする請求項 1 に記載の方法。

40

【請求項 7】

前記繰り返しクラスタリングするステップは、第 1 の反復のクラスタリング結果を前記マルチ型データオブジェクトのすべての関係のあるデータオブジェクトに前記処理ユニットが伝播させるステップを含み、前記関係のあるデータオブジェクトの少なくとも 2 つは異質のデータの型であり、前記クラスタリング結果により形成される特徴空間は、強化されたクラスタリング操作の第 2 の反復において前記マルチ型データオブジェクトのそれぞれのオブジェクトのクラスタリングを高めるために使用されることを特徴とする請求項 1 に記載の方法。

【請求項 8】

前記繰り返しクラスタリングするステップは、前記マルチ型データオブジェクトの個々

50

のオブジェクトの間の類似度を前記処理ユニットが決定するステップを含み、前記類似度はオブジェクト間およびオブジェクト内のコンテンツの類似度、ならびに前記識別された関係の少なくとも1つの間の類似度の少なくとも1つの関数であることを特徴とする請求項1に記載の方法。

【請求項9】

前記繰り返しクラスタリングするステップは、前記マルチ型データオブジェクトの関係のあるオブジェクトをマージして、前記関係のあるオブジェクトの特徴空間の次元を前記処理ユニットが縮小するステップを含むことを特徴とする請求項1に記載の方法。

【請求項10】

オブジェクト型内および異なるオブジェクト型の間の前記マルチ型データオブジェクトの個々のオブジェクト重要度の相互強化を前記処理ユニットが行うステップをさらに含むことを特徴とする請求項1に記載の方法。

【請求項11】

前記送信される語と前記強化されたクラスタの中のオブジェクトの特徴空間とを比較して、前記提案される検索語を前記処理ユニットが識別するステップをさらに備えたことを特徴とする請求項1に記載の方法。

【請求項12】

コンピュータに、

第1の型の少なくとも1つのオブジェクト、および前記第1の型とは異なる第2の型の少なくとも1つのオブジェクトを含む、マルチ型データオブジェクトの間のレイヤ内およびレイヤ間関係の少なくとも1つを決定する手順と、

前記少なくとも1つの関係によって前記マルチ型データオブジェクトを繰り返しクラスタリングすることにより強化されたクラスタ群を生成する手順と、

ユーザから送信される語に関係のある、提案される検索語のリストを、前記強化されたクラスタ群を利用して生成する手順であって、前記提案される検索語は、前記ユーザから送信される語の受信に応答して生成される、手順と、

次の

【数2】

$$\left\{ \begin{array}{l} a(X) = \beta L_X^T h(X) + (1 - \beta) L_{XY} i(Y) \\ h(X) = \beta L_X a(X) + (1 - \beta) L_{XY} i(Y) \\ i(X) = a(X) + h(X) \\ a(Y) = \gamma L_Y^T h(Y) + (1 - \gamma) L_{YX} i(X) \\ h(Y) = \gamma L_Y a(Y) + (1 - \gamma) L_{YX} i(X) \\ i(Y) = a(Y) + h(Y) \end{array} \right.$$

を使用してノードの権威スコアおよびハブスコアを更新することによって、オブジェクト型内および異なるオブジェクト型の間での前記マルチ型データオブジェクトの個々のオブジェクト重要度の相互強化を行う手順であって、 $X = \{x_1, x_2, \dots, x_m\}$ および $Y = \{y_1, y_2, \dots, y_n\}$ は R_X 、 R_Y 、 R_{XY} 、および R_{YX} の関係を有する異質のオブジェクト型のそれぞれのオブジェクトの集合を表現し、方向性が考慮される場合、 L_X および L_Y はそれぞれ集合 X および Y 内の関係を識別するリンクの隣接する行列を表現し、 L_{XY} および L_{YX} は X の中のオブジェクトから Y の中のオブジェクトへの関係を識別するリンクの隣接する行列を表現し、 $a(X)$ および $h(X)$ はそれぞれ X 内のノードの権威スコアおよびハブスコアであり、 $a(Y)$ および $h(Y)$ はそれぞれ Y 内のノードの権威スコアおよびハブスコアを表し、 $i(X)$ および $i(Y)$ はそれぞれ X および Y の中のノードの重要度を表し、およびは異なる関係に由来するリンクの影響を調整するための重みパラメータである、手順と

を実行させるためのプログラムを記録したコンピュータ読取可能な記録媒体。

【請求項 1 3】

前記レイヤ間関係は、コンテンツに関係のある情報、関連トピックに対するユーザの関心、および関連 Web ページに対するユーザの関心のうち 1 つを含むことを特徴とする請求項 1 2 に記載のコンピュータ読取可能な記録媒体。

【請求項 1 4】

前記レイヤ内関係は、クエリ改良、推奨される Web ページ、およびそれぞれのユーザの間の関係のうち少なくとも 1 つを含むことを特徴とする請求項 1 2 に記載のコンピュータ読取可能な記録媒体。

【請求項 1 5】

前記マルチ型データオブジェクトのそれぞれは、検索クエリデータオブジェクト型、選択された Web ページ型、およびユーザ情報型のうち少なくとも 1 つに関係のあることを特徴とする請求項 1 2 に記載のコンピュータ読取可能な記録媒体。

10

【請求項 1 6】

前記決定された関係のうち少なくとも 1 つに重みを付けることによって、前記マルチ型データオブジェクトの関連するオブジェクトに対する重要度を示すことを特徴とする請求項 1 2 に記載のコンピュータ読取可能な記録媒体。

【請求項 1 7】

前記識別する手順および繰り返しクラスタリングする手順は、検索語提案のために実行されることを特徴とする請求項 1 2 に記載のコンピュータ読取可能な記録媒体。

【請求項 1 8】

前記繰り返しクラスタリングする手順は、第 1 の反復のクラスタリング結果を前記マルチ型データオブジェクトのすべての関係のあるデータオブジェクトに伝播させる手順を含み、前記関係のあるデータオブジェクトの少なくとも 2 つは異質のデータの型であり、前記クラスタリング結果により形成される特徴空間は、強化されたクラスタリング操作の第 2 の反復において前記マルチ型データオブジェクトのそれぞれのオブジェクトのクラスタリングを高めるために使用されることを特徴とする請求項 1 2 に記載のコンピュータ読取可能な記録媒体。

20

【請求項 1 9】

前記繰り返しクラスタリングする手順は、前記マルチ型データオブジェクトの個々のオブジェクトの間の類似度を決定する手順を含み、前記類似度は、オブジェクトコンテンツの類似度および前記識別された関係の少なくとも 1 つの間の類似度の少なくとも 1 つの関数であることを特徴とする請求項 1 2 に記載のコンピュータ読取可能な記録媒体。

30

【請求項 2 0】

前記繰り返しクラスタリングする手順は、前記マルチ型データオブジェクトの関係のあるオブジェクトをマージすることによって前記関係のあるオブジェクトの特徴空間の次元を縮小する手順を含むことを特徴とする請求項 1 2 に記載のコンピュータ読取可能な記録媒体。

【請求項 2 1】

前記コンピュータに、オブジェクト型内および異なるオブジェクト型の間で前記マルチ型データオブジェクトの個々のオブジェクトの重要度の相互強化を行う手順をさらに実行させることを特徴とする請求項 1 2 に記載のコンピュータ読取可能な記録媒体。

40

【請求項 2 2】

前記コンピュータに、前記送信される語と前記強化されたクラスタの中のオブジェクトの特徴空間とを比較して、前記提案される検索語を識別する手順を実行させることを特徴とする請求項 1 2 に記載のコンピュータ読取可能な媒体。

【発明の詳細な説明】

【技術分野】

【0001】

本発明は、データマイニングに関し、より詳細には、異質のオブジェクトをクラスタリングすることにより、検索語提案のためのシステムおよび方法を高めることに関する。

50

【背景技術】

【0002】

キーワードまたは句は、WWW (World Wide Web) 上で関連するWeb ページ/サイトを検索する場合、Webサーファ (surfer) によって検索エンジンに送信される単語または語の集合である。検索エンジンは、ページ/サイト上に現れるキーワードおよびキーワード句に基づき、Webサイトの関連性を決定する。かなりの割合のWebサイトトラフィックは、検索エンジンの使用に起因するので、Webサイト主催者は、適切なキーワード/句選択が、所望のサイトの登場を獲得するためにサイトトラフィックの増加に不可欠であることを理解している。検索エンジン結果を最適化するためにWebサイトと関係のあるキーワードを識別する技術には、例えば、関係のあるキーワードを識別するためにWebサイトのコンテンツおよび目的の人手による評価を含む。この評価には、キーワード人気度ツール (keyword popularity tool) の使用を含むことができる。そのようなツールは、どれだけの数の人々が、特定のキーワードまたは特定のキーワードを含む句を検索エンジンに送信したかを決定する。Webサイトと関係があり、検索クエリを生成する際により頻繁に使用されると決定されたキーワードを、一般に、Webサイトに関して、検索エンジン結果を最適化するために選択する。

10

【0003】

Webサイトの検索エンジン結果を最適化するためのキーワードの集合を識別した後、主催者は、(他のWebサイトの検索エンジン結果の表示位置と比較して) 検索エンジン結果においてより上位にWebサイトを昇進させることを所望する可能性がある。この目的のために、主催者は、(複数の) キーワードを入札することにより、(複数の) キーワードに結びつく主催者のリストをWebサーファがクリックするたび毎に主催者がいくら支払うかを表す。すなわち、キーワード入札は、ペイパークリック (pay-per-click) 入札である。同一のキーワードに対する他の入札と比較してキーワード入札の総額がより大きいほど、検索エンジンにより、キーワードに基づく検索結果において結びつくWebサイトがより上位に(重要度に関してより目立つように) 表示される。

20

【0004】

Webサイトのコンテンツと関係があり、入札する(複数の) 語を識別するための従来システムおよび技術は、通常、クラスタリングアルゴリズムを使用することにより、同じクラスタからのオブジェクトは類似し、異なるクラスタからのオブジェクトは類似しないような仕方で、オブジェクトの集合をグループまたはクラスタに分割する。クラスタリングアプローチでは、クラスタリングされるデータオブジェクトが、独立であり、同一のクラスのオブジェクトであり、しばしば、固定長ベクトルの特徴/属性値でモデル化されていると仮定する。最近急増しているデータマイニング研究において、この古典的な問題が、大型データベースの文脈において再検討されている。しかし、クラスタリングされるデータオブジェクトの均質性は、たとえ、Webマイニングおよび協調フィルタリングなどいくつかのアプリケーションの出現により、そのような仮定に対する異議を唱えていても、依然、基本的な仮定であるように思われる。そのようなアプリケーションでは、データオブジェクトは、異なる型のデータオブジェクトであり、非常に相互に関係がある。残念ながら、たとえ異質のオブジェクト型にわたって分けられたオブジェクトが非常に相互に関係がある可能性があっても、従来のクラスタリング操作は、通常、それぞれのオブジェクト型を個々に、異なるオブジェクト型の相互に関係がある様相を考慮することなしにクラスタリングする。

30

40

【0005】

このことの1つの理由は、異なる型のデータオブジェクト間の関係が、希薄であり識別するのが困難である場合が多いからである。別の理由は、それぞれのオブジェクトに結びついた静的な固定長の値のベクトルを有する、あらゆるそのような関係の表現は、ここで、ベクトルはオブジェクト属性および異なる型の関係のあるオブジェクトの属性の両方を表すが、非常に高次元(特徴空間)のオブジェクト属性/特徴ベクトルを作り出すからで

50

ある。このような高次元は望ましくない。なぜならば、特徴空間内でデータが互いに遠く離れ、効率的なモデルを小さい領域内の希薄な量のデータで十分に扱うことができないからである。

【発明の開示】

【発明が解決しようとする課題】

【0006】

したがって、異質のデータオブジェクトにわたる関係の観点から、関係のあるオブジェクト（例えば、語）を識別しグループ化する、より良いクラスタリング技術は有用である。クラスタリング技術を使用することにより、例えば、検索エンジン最適化および語の入札のための（複数の）語を識別するシステムおよび方法を提供し、それによってシステムと方法の両方において、関係のある（複数の）語を識別する確率を大いに高めることができる。

10

【課題を解決するための手段】

【0007】

関係のある語提案のためのシステムおよび方法を説明する。一態様では、2つ以上のマルチ型データオブジェクトのうちのそれぞれのオブジェクトの間にレイヤ内および/またはレイヤ間関係を識別する。マルチ型データオブジェクトのそれぞれのオブジェクトには、第1型の少なくとも1つのオブジェクト、および第1型とは異なる第2型の少なくとも1つのオブジェクトを含む。マルチ型データオブジェクトを、関係のうちのそれぞれの関係の観点から繰り返しクラスタリングすることによって、強化されたクラスタを生成する。

20

【0008】

図では、構成要素参照符号の左端の数字が、その構成要素が最初に現れる特定の図を識別する。

【発明を実施するための最良の形態】

【0009】

（概要）

図1は、相互に関係のある異質のオブジェクトデータ型の例示的フレームワーク100を示す。フレームワーク100は、異質のデータオブジェクト/ノードの複数のレイヤ102、ならびに結びつけられたレイヤ間およびレイヤ内データオブジェクトリンク/関係を含む。各レイヤ102-1ないし102-Nは、同一型（均質の）データオブジェクトまたはノードのそれぞれの集合を含む。すなわち、ノード集合Pは、同一データ型のおのである1つまたは複数のデータオブジェクト p_1 ないし p_j を含み、ノード集合Uは、同一データ型のおのである1つまたは複数のデータオブジェクト u_1 ないし u_k を含み、以下同様である。このため、異なるそれぞれのレイヤ102にあるデータオブジェクトの型は、お互いに関して異質である。

30

【0010】

この実装では、例えば、

- ・レイヤ102-1は、（複数の）検索クエリデータオブジェクト/（複数の）ノード p_1 ないし p_j を含むマイニングされた（mined）検索クエリである。検索クエリオブジェクトは、（複数の）クエリの語を含み、以下に説明するとおり、クエリログからマイニングされた複数の履歴クエリのそれぞれのクエリを表す。
- ・レイヤ102-2は、（複数の）Webページデータオブジェクト/（複数の）ノード u_1 ないし u_k を含むマイニングされたWebページレイヤである。
- ・レイヤ102-3は、マイニングされたユーザレイヤであり、（複数の）ユーザ情報オブジェクト/複数のノード w_1 ないし w_m を含む。
- ・レイヤ102-Nは、それぞれの異なるオブジェクト型 x_1 ないし x_o を備えるレイヤ102はあらゆる個数が可能であることの例示を示す。

40

【0011】

一組のデータオブジェクトの間に張るライン/リンクは、それぞれのデータオブジェク

50

トの間に存在すると決定されたそれぞれのマイニングされた関係を表す。クラスタリングの所定の実施形態では、ライン/リンクを「エッジ」と呼ぶ。本明細書では、一般化された語のラインまたはリンクを使用することにより、リンク、エッジ、またはオブジェクト間の関係を説明するあるオブジェクトと別のオブジェクトとのあらゆる連結子を説明する。リンク方向を（それは、データオブジェクト間の関係を示す矢印により提供されるが）、関与する（participating）オブジェクト属性の相関的要素として、いずれかの方向に向けることができる。リンクは、例示的であり、範囲を限定するものではないと見なされる。フレームワーク100で表されるようなWeb環境における所定のリンクを、1つの方向により適切に向けることができ、矢印の方向は、通常、以下に説明する強化されたクラスタリング操作に影響を与えない。

10

【0012】

オブジェクトの組の間のリンクを、レイヤ内またはレイヤ間リンクとして分類することができる。レイヤ内リンクは、同一型の異なるオブジェクトの間の識別された関係を例示する。このように、レイヤ内リンク104は、同一レイヤ102内のオブジェクトを接続する。例えば、データオブジェクトのそれぞれの組の間における実線104は、レイヤ内リンクを表す。この例では、レイヤ内リンクは、Webページオブジェクト u_2 から別のWebページオブジェクト u_3 に張り、異なるWebページ間の（複数の）関係を表す。

【0013】

レイヤ間リンクは、異なる型のデータオブジェクト間の関係を説明する。レイヤ間リンクは、異質のオブジェクトから成る組のそれぞれのオブジェクト間に張られるため、データオブジェクトから成る関与する組のおおのを、異なるそれぞれのデータオブジェクト/ノード集合レイヤ102上に示す。図1に示すとおり、実線ではないオブジェクトからなる組を接続するあらゆるラインは、レイヤ間リンクである。例えば、リンク106は、オブジェクトから成る組の第1のオブジェクトから、オブジェクトから成る組の第2のオブジェクトへの参照（例えば、ハイパーリンク）を示し、リンク/ライン108は、オブジェクトから成る組の第1のオブジェクトから、オブジェクトから成る組の第2のオブジェクトに共有される/参照される問題（例えば、主題要素）を示し、リンク/ライン110は、オブジェクトから成る組の第1のオブジェクトから、オブジェクトから成る組の第2のオブジェクトへのブラウズリンクを示す。別の例では、リンクを、ユーザオブジェクト w_4 から検索クエリオブジェクト p_5 およびWebページオブジェクト u_5 へ張り、ユーザにより関係があるとして選択されたWebページを戻すクエリをユーザが送信することを表すことが可能である。

20

30

【0014】

図1の実施例では、レイヤ内およびレイヤ間リンクのそれぞれのリンクで示されるとおり、異なるオブジェクト型（p、u、w、...）が関係している。例えば、ユーザ（オブジェクトwで表される）が、クエリ（オブジェクトp）を発行する。ユーザは、発行されたクエリの受信に回答して検索エンジンによって戻されたWebページ（オブジェクトu）をブラウズし、各検索クエリ（オブジェクトp）は、1つまたは複数のそれぞれのWebページ（オブジェクトu）を参照する。以上の観点から、Webユーザ情報をクラスタリングする場合、ユーザがブラウズした（複数の）Webページ、およびそれぞれの（複数の）Webページを獲得するのに使用されたクエリは、より高い類似性を有し、クラスタリングプロセスにおいて一緒にクラスタリングされる傾向にあるはずである。同様に、Webページをクラスタリングする場合、Webページがどのようにユーザによって使用され、どのようにそれぞれの検索クエリによって参照されたかも考慮に入れなければならない。これに対処するため、以下に説明するとおり、強化されたクラスタリングアルゴリズムが、複数のデータオブジェクトのそれぞれのオブジェクトの間のマイニングされた関係の相関的要素として、そのような異質のデータオブジェクトをクラスタリングする。

40

【0015】

本発明の一態様は、本来備わっている相互関係に基づき、クラスタリングされるオブジェクトに、他のオブジェクトへのリンクを提供する。各オブジェクトに接続する複数のリ

50

ンク（およびそれらのリンクが接続する先の複数のオブジェクト）の所定のリンクには、そのオブジェクトに対するリンクの関連性を反映するように異なる重要度で重みを付けることができる。例えば、クラスタリングされた複数のオブジェクトと同一型の複数のオブジェクトには、異なる型の複数のオブジェクトより大きい重要度を提供することができる。本発明は、様々なレベルの重要度を異なるオブジェクトまたは異なる型のオブジェクトに割り当てることができることによる仕組みを提供する。異なるオブジェクト（または異なる型のオブジェクト）に異なるレベルの重要度を割り当てて、本明細書では、重要度付きクラスタリングと呼ぶ。異なるオブジェクトの様々なレベルの重要度により、クラスタリングの結果および効果を改善する場合が多い。次に、検索語提案のためのマルチ型データオブジェクトの強化されたクラスタリングの、以上および他の態様を説明する。

10

【0016】

Webサイトと関係があり、検索クエリを生成する際にエンドユーザによってより頻繁に使用されたと判定された（複数の）語／（複数の）キーワードは、一般に、Webサイトに関して検索エンジン結果を最適化するために、Webサイト主催者／広告主によって選択される。これに留意して、以下の本発明のシステムおよび方法は、本実装において、検索語提案であり手近なタスクと相互に関係があると決定されたマルチ型データオブジェクトをマイニングする。そのようなマルチ型データオブジェクトには、履歴クエリを検索エンジンに送信することによって獲得された結果からマイニングされた意味論的文脈（例えば、テキスト、URL、結果タイトル、および各結果の短い説明など）で高められマイニングされた履歴検索クエリの（複数の）語、特定の履歴検索クエリに回答してユーザによって選択されたWebページの集合、ユーザに固有の情報（例えば、ユーザのWebサイトアクセス情報、検索クエリを生成するのに使用されたマシンのIPアドレスなど）、および／または関係のあるデータオブジェクトの型を含む。

20

【0017】

マルチ型データオブジェクト間の類似度を、識別され重みが付けられたコンテキストの類似度と、計算されたオブジェクト間およびオブジェクト内関係の類似度との線形結合として決定する。データオブジェクトには、オブジェクト間および型内関係に由来するリンク構造を分析することにより、異なる重みを割り当てる。したがって、複数のデータオブジェクトのそれぞれのオブジェクト間の類似度には、オブジェクト自体の属性の類似度だけでなく、オブジェクトの関係の類似度も含む。

30

【0018】

以上の計算されたマルチ型オブジェクト関係の観点から、強化されたクラスタリングアルゴリズムは、各オブジェクトの識別されたオブジェクト間およびオブジェクト内関係属性の相関的要素として、マルチ型データオブジェクトを繰り返しクラスタリングする。本実装において、クラスタ内のオブジェクトの重み付けされた合計を使用することによってクラスタ中心（centroid）を算出するのに、変更された直接k平均（direct-k-means）アルゴリズムを使用する。これは、オブジェクトのそれぞれの関係属性を更新することによってすべての関係のあるデータオブジェクトにクラスタリング結果を伝播させる反復プロセスである。すなわち、1つの型のオブジェクトのクラスタリング結果が、新たな特徴空間を形成し、次に、この空間を、他の関係のある異なる型のオブジェクトに射影し、伝播させる。次に、関係のある型のオブジェクトに対するクラスタリングを、その更新された特徴空間を使用して実行する。この反復強化プロセスを、各オブジェクト型に対して実行することにより、特徴空間の次元を減らすよう大いに関係のあるクラスタノードをマージし、すべてのマルチ型オブジェクトにわたるクラスタリング結果が収束するまで続ける。これにより、大いに高い関係のあるマルチ型データオブジェクトの強化されたクラスタとなる。

40

【0019】

エンドユーザから語を受け取ることに回答して、システムおよび方法は、語／クエリオブジェクト型に基づき、強化されたクラスタ内の（複数の）語のそれぞれの語とその（複数の）語を比較する。強化された語のクラスタは、互いに文脈上、別の語と関係のある語

50

を含むので、送信された入札をクラスタ内の語と比較する場合、語句を複数の関係のある文脈、すなわち「意味」の観点から評価する。さらに、各々の強化された語のクラスタは、高い関係のあるマルチ型オブジェクトの集合に由来するので、アルゴリズムは、純粋に文脈ベースの方法の欠点を克服することができる。すなわち、クエリの語の間の意味関係を効率的に強化し、語の文脈における雑音の影響を抑えることができる。受け取られた語と強化されたクラスタ内のオブジェクトの特徴空間とを比較することに対応して、1つまたは複数の検索語提案を識別する。これらの検索語提案をエンドユーザに伝達する。

【0020】

(例示的システム)

必要ではないが、本発明は、パーソナルコンピュータによって実行されるコンピュータ実行可能命令(プログラムモジュール)の一般的な文脈で説明する。プログラムモジュールには、一般に、特定のタスクを実行する、または特定の抽象データ型を実装するルーチン、プログラム、オブジェクト、コンポーネント、データ構造などが含まれる。本システムおよび方法を上述の文脈で説明するが、以下に説明する動作および操作は、ハードウェアに実装することもできる。

【0021】

図2は、検索語提案のためのマルチ型データオブジェクトの強化されたクラスタリングを行うための例示的システム200を示す。本実装において、システム200は、ネットワーク204を介してクライアントコンピュータ処理装置206に結合された編集確認サーバ(EVS)202を含む。例えば、クライアントコンピュータ処理装置206またはEVS202上で実行される別のアプリケーション(図示せず)から、(複数の)語208を受け取ったことに対応して、EVS202は、提案される語リスト210を生成し、クライアントコンピュータ処理装置206に伝達することにより、エンドユーザが、実際に(複数の)語の入札を行う前に、(複数の)語208に意味的および/または文脈的に関係のある語の集合を評価することができる。ネットワーク204は、オフィス、企業全体のコンピュータネットワーク、イントラネット、およびインターネットで一般的であるような、ローカルエリアネットワーク(LAN)通信環境と一般的なワイドエリアネットワーク(WAN)通信環境とのあらゆる組み合わせを含むことが可能である。システム200がクライアントコンピュータ処理装置206を含む場合、クライアントコンピュータ処理装置は、パーソナルコンピュータ、ラップトップ、サーバ、モバイルコンピュータ処理装置(例えば、セルラー電話機、パーソナルデジタルアシスタント、またはハンドヘルドコンピュータ)などのあらゆる型のコンピュータ処理装置である。

【0022】

EVS202は、提案される語リスト210を生成するいくつかのコンピュータプログラムモジュールを含む。複数のコンピュータプログラムモジュールには、例えば、検索語提案(STS)モジュール212が含まれる。本実装において、説明および例示的図解の目的で、STSモジュール212は、履歴クエリの語マイニング、Webページ検索、特徴抽出、特徴空間次元の縮小および正規化、マルチ型データオブジェクトの強化されたクラスタリング、ユーザ入札の(複数の)語を強化されたクラスタの内容とマッチングすることにより検索語提案を実行すること、および語分類などの、複数の機能を実行することとして説明する。以上の複数の操作のそれぞれの操作は、STSモジュール212と通信する1つまたは複数の他のコンピュータプログラムモジュール(図示せず)によって実行されることも可能であることを認めることができよう。

【0023】

(意味論的文脈でマイニングされた履歴検索クエリの強化)

STSモジュール212は、本実装において、検索語提案である、手近なタスクと相互に関係があると決定された複数のマルチ型データオブジェクト(MDO)214をマイニングする。そのような複数のマルチ型データオブジェクト214には、複数の履歴クエリ216のそれぞれのクエリを検索エンジンに送信することによって獲得された検索結果からマイニングされた意味論的文脈(例えば、テキスト、URL、結果タイトル、および各

結果の短い説明など)でSTSモジュール212によって高められ、マイニングされた複数の履歴検索クエリ216の(複数の)語、および特定の履歴検索クエリに回答してユーザによって選択されたWebページの集合、ユーザに固有の情報(例えば、ユーザのWebサイトアクセス情報、検索クエリを生成するのに使用されたマシンのIPアドレスなど)、および/または関係のあるマルチ型データオブジェクトの型などの「他のMDO」218を含む。

【0024】

詳細には、STSモジュール212は、(複数の)クエリログ220から履歴クエリ216の集合を検索する。(複数の)履歴クエリ216は、1名または複数名のユーザによって検索エンジンに以前に送信された検索クエリの語を含む。STSモジュール212は、出現頻度の相関的要素として履歴クエリ群216を評価することより、高いFOO(出現頻度)の検索語222と比較的低い出現頻度の検索語224を識別する。本実装において、設定可能な閾値を使用することにより、履歴クエリが比較的高い出現頻度を有するか、または比較的低い出現頻度を有するかを決定する。例えば、少なくとも閾値回数、出現する複数の履歴クエリ216の中の検索クエリの語は、高い出現頻度を有するとされる。類似して、閾値回数より少ない回数、出現する複数の履歴クエリ216の中の検索クエリの語は、低い出現頻度を有するとされる。例示の目的で、そのような閾値を「他のデータ」226のそれぞれの部分として示す。

【0025】

STSモジュール212は、各クエリを1つずつ(検索クエリ227)、検索エンジン228に送信することにより、高い出現頻度のクエリの語222の意味論上/文脈上の意味をマイニングする。検索クエリ227を受け取ったことに回答して、検索エンジン228は、(複数の)検索結果230の中でランク付けされたリスト(数が設定可能である)をSTSモジュール212に戻す。ランク付けされたリストは、URL、結果タイトル、ならびに送信された検索クエリ227に関係のあるクエリの語の短い説明および/または文脈を含む。ランク付けされたリストは、検索結果230の中に格納される。そのような検索結果の検索が、各検索クエリ227に関して行われる。

【0026】

STSモジュール212は、WebページHTML(ハイパーテキストマークアップ言語)を解析して、検索された各検索結果230から、各クエリの語222に関するクエリの語のURL、結果タイトル、ならびに短い説明および/または文脈を抽出する。クエリの語のURL、結果タイトル、短い説明および/または文脈、ならびに検索された検索結果230を獲得するのに使用された検索クエリ227は、抽出された特徴232のそれぞれのレコードの中にSTSモジュール212によって格納される。

【0027】

高い出現頻度のクエリの語222に関して検索結果230を解析した後、STSモジュール212は、抽出された特徴232に対してテキスト前処理操作を実行して、抽出された特徴から個々のキーワードに入る言語トークンを生成する(トークン化する(tokenize))。トークンの次元を縮小するため、STSモジュール212は、例えば、ポーターステミング(Porter stemming)アルゴリズムを使用して、あらゆるストップワード(stop-word)(例えば、「the」、「a」、「is」など)を削除し、一般的な接尾辞を削除することにより、キーワードを正規化する。STSモジュール212は、もたらされる抽出された特徴232を、1つまたは複数の語に基づくMDO(マルチ型データオブジェクト)ベクトル234に構成する。

【0028】

各語に基づくマルチ型データオブジェクトベクトル234は、語の頻度に基づく次元、および逆ドキュメント頻度(TFIDF)スコアを有する。第i番のベクトルの第j番のキーワードに対する重みは、以下のとおり計算される。すなわち、

$$w_{ij} = TF_{ij} \times \log(N/DF_j)$$

ただし、 TF_{ij} は、語の頻度(第i番のレコード内のキーワードjの出現回数)を表し

10

20

30

40

50

、 N は、クエリの語の総数であり、 DF_j は、キーワード j を含むレコードの数である。

【0029】

各クエリの語のベクトル表現を所与として、コサイン関数を使用することにより、一組の語の間の類似度を測定する（ベクトルは正規化されていることを想起されたい）。すなわち、

【0030】

【数1】

$$\text{sim}(q_j, q_k) = \sum_{i=1}^d w_{ij} \cdot w_{ik}$$

10

【0031】

すなわち、2つの語の間の距離（類似測度）は、以下のとおり定義される。すなわち、

$$\text{dist}(q_j, q_k) = 1 - \text{sim}(q_j, q_k)$$

このような類似測度を、「他のデータ」226のそれぞれの部分として示す。例示的なこのような類似度値を、以下に説明する表1の例示的な提案される語のリスト210の中で示す。

【0032】

（ユーザが選択したWebページおよびユーザ情報のマイニング）

Webサイトの検索エンジン結果を最適化するために実質的に最も関係のある語の集合（検索語提案）を識別するため、STSモジュール212は、複数の履歴クエリ216とは異なる／異質のマルチ型データオブジェクト214をマイニングする。説明の目的のために、これらのマイニングされる複数のオブジェクトを「他のMDO」218と表す。所定の実装において、「他のMDO」218には、例えば、エンドユーザが選択したWebページおよび／またはユーザ固有の情報が含まれ、ただし、ユーザは、履歴クエリ216を検索エンジン228に送信することに関連するユーザである。STSモジュール212は、クエリログ220から、エンドユーザが選択したWebページを抽出する。エンドユーザが選択したWebページは、希薄であることも、そうでないことも可能であり、希薄は平均すると、例えば、履歴クエリ216当たり2から3のWebページになる。STSモジュール212は、（複数の）クエリログ220または他のデータソースからユーザ固有の情報を抽出する。ユーザ固有の情報には、例えば、複数の履歴クエリ216のそれぞれのクエリを送信するのに使用されたマシンのIP（インターネットプロトコル）アドレス、GUID、および／またはWebサイトアクセス情報（例えば、マイクロソフト社のドットネットパスポートの情報）が含まれる。

20

30

【0033】

（強化されたマルチ型データオブジェクトのクラスタリング）

STSモジュール212は、クラスタリング分析のためにマルチ型の相互に関係のあるデータオブジェクト（MDO214）間の関係を十分に詮索する。複数のマルチ型データオブジェクト214には、 n 個の異なる型のオブジェクト、 X_1 、 X_2 、 \dots 、 X_n （例えば、複数の履歴クエリ216および「他のMDO」218）が含まれる。各型のデータオブジェクト X_i は、特徴集合 F_i で記述される。同一型内の複数のデータオブジェクトは、型内関係、 $R_i: X_i \times X_i$ で相互に関係がある。2つの異なる型からのデータオブジェクトは、型間関係、 $R_{ij}: X_i \times X_j$ を介して関係がある。これらの関係と区別するため、 F_i を、データオブジェクトのコンテンツ特徴（content feature）と呼ぶ。特定のオブジェクト $x \in X_i$ に関して、 $x \in F_i$ を使用してそのオブジェクトのコンテンツ特徴を表し、 $x \in R_i: X_i$ および $x \in R_{ij}: X_j$ を使用して、それぞれ X_i および X_j の中でそのオブジェクトに関係のあるオブジェクトを表す。マルチ型の相互に関係のあるデータオブジェクトをクラスタリングすることの問題は、各型の複数のオブジェクト X_i を K_i 個のクラスタに分けて、各クラスタ内のデータオブジェクトが高い類似度を有し、異なるクラスタからのオブジェクトが似ていないようにすることであ

40

50

る。

【 0 0 3 4 】

マルチ型データオブジェクト群 2 1 4 のオブジェクトが、コンテンツ特徴と、複数のマルチ型データオブジェクト 2 1 4 の中の他の (複数の) オブジェクトとの関係をとともに有することを考慮すると、2つのオブジェクト間の類似度を、以下の数式に従って決定する。すなわち、

$$S = \alpha \cdot S_f + \beta \cdot S_{i n t r a} + \gamma \cdot S_{i n t e r} \quad (1)$$

ただし、 S_f は、コンテンツ類似度であり、 $S_{i n t r a}$ および $S_{i n t e r}$ はそれぞれ、型内類似度および型間類似度であり、 α 、 β 、および γ は、異なる類似度に対する重みであり、 $\alpha + \beta + \gamma = 1$ である。

10

【 0 0 3 5 】

(1) から、2つのオブジェクト間の類似度は、コンテンツ類似度と関係類似度の線形結合である。 α 、 β 、および γ に異なる値を割り当てることにより、S T S モジュール 2 1 2 は、全体的な類似度における異なる類似度の重みを調整 / 構成することができる。例えば、 $\alpha = 1$ 、 $\beta = \gamma = 0$ である場合、コンテンツ特徴間の類似度が考慮される。 $\alpha = 0$ に設定することにより、S T S モジュール 2 1 2 は、型内類似度の効果を顧慮しない。

【 0 0 3 6 】

等式 1 における類似度は、通常、オブジェクトの型およびアプリケーションによって決定され、異なる関数を使用して定義することができる。例えば、2つの Web ページ間のコンテンツ類似度は、Web ページのコンテンツに由来する 2つのキーワードベクトルの

20

【 0 0 3 7 】

特定のオブジェクトの関係特徴は、エントリが特定のオブジェクトに関係のあるオブジェクトに対応する M D O ベクトル 2 3 4 によって表される。所定の実装において、各エントリは、関係の重みに対応する数値である。例えば、2つのオブジェクト型、 $X = \{ x_1, x_2, \dots, x_m \}$ および $Y = \{ y_1, y_2, \dots, y_n \}$ を所与として、オブジェクトの型間関係は、 $V_X = [v_1, v_2, \dots, v_n]^T$ として定義され、ただし、 $v_i \geq 0$ の場合であり、それ以外の場合、 $v_i = 0$ である。すると、 X の中の 2つのオブジェクト間の型間関係 $R_{X Y}$ 上の類似度 $S_{i n t e r - X Y}$ も、2つのベクトルのコサイン関数として定義することが可能である。

30

【 0 0 3 8 】

X_i の中のオブジェクトが、複数のデータオブジェクト型と型間関係を有する場合、最終的な型間類似度は、すべての型間類似度の線形結合であることが可能である。

【 0 0 3 9 】

定義された類似度関数を使用して、S T S モジュール 2 1 2 は、複数の履歴クエリ 2 1 6 および「他の M D O 」 2 1 8 の間のレイヤ内関係 / リンクおよびレイヤ間リンクを識別する。クラスタリングにおけるレイヤ間リンクの使用は、所定の型のオブジェクトのクラスタリングが別の型のオブジェクトによって影響を及ぼされる可能性があることを認識する。例えば、Web ページオブジェクトのクラスタリングは、ユーザオブジェクトの構成、状態、および特性によって影響を及ぼされる可能性がある。したがって、それらのマイ

40

ニングされたレイヤ内関係およびレイヤ間関係は、以下に説明するとおり、相互に関係のあるデータオブジェクトのクラスタ品質を向上させるのに使用される。マイニングされたレイヤ間およびレイヤ内データオブジェクト関係は、各オブジェクトのそれぞれの M D O ベクトル 2 3 4 の中に格納される。

【 0 0 4 0 】

所定の実装において、識別されたレイヤ間リンク / 関係は、例えば、以下の 1 つまたは複数を表す。すなわち、

- ・コンテンツに関係のある情報、例えば、履歴クエリ 2 1 6 内のリンク、およびユーザが選択した (クリックスルーされた (c l i c k e d - t h r u)) Web ページに対応するリンク

50

・例えば、履歴クエリ 2 1 6 内のリンク、およびユーザ固有の情報により決定された、関連するトピックに対するユーザの関心

・例えば、ユーザ固有の情報と選択された Web ページの間のリンクを介して決定された、選択された Web ページに対するユーザの関心

所定の実装において、識別にされたレイヤ内リンク / 関係 (同一のデータ型のオブジェクト間の関係) は、例えば、以下の 1 つまたは複数を表す。すなわち、

・以下により詳細に説明する、クエリ内のリンク

・ユーザが選択した Web ページ内の内向き / 外向きの (directed in / out) ハイパーリンクで示される、推奨される (複数の) Web ページ

・例えば、それぞれのユーザの間で識別された関係 / リンクで示される人間関係。所定の実装において、この型の関係情報は、ユーザプロファイルの計算された類似度からマイニングされる。ユーザプロファイルには、例えば、人口統計、地理的位置、 (複数の) 関心などが含まれる。所定の実装において、ユーザプロファイルは、それぞれのユーザが供給する Web サイトアクセス情報を介してアクセスされる。

【 0 0 4 1 】

クエリ内のリンクに関して、クエリ内のリンクによって示されるレイヤ内関係は、初期の履歴クエリ 2 1 6 と後続のクエリ改良 (refinement) (複数の履歴クエリ 2 1 6 のそれぞれのクエリによっても表現される) との間の、または後続のクエリ改良間のリンクを表現する。所定の実装において、そのような情報は、 (複数の) クエリログ 2 2 0 から検索されたクリックスルー Web ページ情報から抽出される。より詳細には、初期検索クエリ結果が満足 of いくものでないと決定されると、ユーザは、初期クエリが送信された時点から設定可能な合計時間内に、 1 つまたは複数の改良されたクエリを検索エンジン 2 2 8 へ送信することが想定される。設定可能な合計時間は、クエリセッションを表現する。 1 回または複数回のそのような検索クエリの語の改良の後、ユーザは、満足 of いく検索結果を獲得することができる。例えば、ユーザが、製品サポートの Web サイトを訪れ、 「クッキー (cookie) 」 という初期クエリを送信することを考慮されたい。検索結果が満足 of いくものでない (例えば、広すぎる) と決定されると、ユーザは、 「クッキーを可能にする (enable cookie) 」 クエリの語に変更 / 改良して、より満足 of いく検索結果を獲得することができる。

【 0 0 4 2 】

所定の実装において、STS モジュール 2 1 2 は、 (複数の) クエリログ 2 2 0 の 1 つまたは複数の部分をそれぞれのクエリセッションにセグメント化することにより、クエリ内のリンクを識別する。各クエリセッションは、初期クエリ、 1 つまたは複数のクエリ改良、およびできる限り 1 つまたは複数の Web ページクリックスルーの指示を含むことが可能である。初期クエリ、および 1 つまたは複数の関連するクエリ改良を分類するため、STS モジュール 2 1 2 は、各クエリセッションのクエリ間の語の類似度を計算する。類似度の 1 つまたは複数の閾値基準を満たす検索クエリが、クエリ内、および対応するクエリ改良内のリンクを生成するために選択される。所定の実装において、クエリの類似度を、例えば、段落 (0 0 3 8) で上述した例示的操作を使用して決定する。

【 0 0 4 3 】

複数のマルチ型データオブジェクト 2 1 4 の間の関係を複数の M D O ベクトル 2 3 4 の対応するベクトルにおいてモデル化された関係特徴としてマッピングした後、従来のクラスタリング技術を使用して (すなわち、本明細書で開示する強化されたクラスタリング操作を使用せずに) 、各型のデータオブジェクトを個々にクラスタリングすることも可能である。しかし、データオブジェクトを個々にクラスタリングすることさえ、最初は、実行可能であるように思われる可能性があるものの、この技術は、実質的に限られており、問題を抱えている。そのことの 1 つの理由は、関係のための特徴ベクトルのサイズが非常に大きくなるにつれ、オブジェクトの数が非常に多くなるためである。また、関係のあるオブジェクトの正確なマッチングに基づく、関係特徴に関して定義された類似度が、 0 でないエントリの希薄さにより悪化する。別の理由は、データオブジェクト間の (複数の) 関

係が、データオブジェクトに割り当てられた特徴に十分に反映されていない可能性があるが、クラスタリングプロセス自体の過程にだけ発見される可能性があることを、従来のクラスタリング技術が考慮しないためである。すなわち、既存のクラスタリング技術は、順にクラスタリング操作を行うことにより、後続の分析／クラスタリング操作においてデータを強化する際に有用な構造化された情報を提供できることを考慮しない。

【 0 0 4 4 】

S T S モジュール 2 1 2 は、1つのデータオブジェクト型のクラスタリング結果をそのオブジェクトと関係のあるすべてのデータオブジェクト型に、それらのデータオブジェクト型のそれぞれの関係特徴を更新することによって少なくとも伝播させることにより、従来のクラスタリング技術の以上の問題／限界に対処する。すなわち、S T S モジュール 2 1 2 は、強化された複数のクラスタ 2 3 6 のコンテンツに基づく個々のマルチ型データオブジェクト 2 1 4 に対する指示されたデータオブジェクトの（複数の）関係を集約することにより、強化された複数のクラスタ 2 3 6 を生成する。例えば、クラスタリング後に2つの候補ノードが存在する場合、最も近接した2つの候補ノードを、例えば、その2つの候補ノードのベクトル値を平均することにより、マージすることができる。このマージが個々のノードを結合することにより、考慮すべきノードの数を減らすことが可能になる。真の意味において、（複数の）M D O ベクトル 2 3 4 の次元を縮小する。次に、S T S モジュール 2 1 2 は、（複数の）M D O ベクトル 2 3 4 をクラスタリングする。このプロセスは、すべてのオブジェクト型のクラスタリング結果が収束するまで、繰り返し実行される。

【 0 0 4 5 】

反復クラスタリング射影技術は、均質の型のオブジェクトを各レイヤが含む、別々のレイヤに構成された別々の型のオブジェクトからクラスタリング情報を獲得することに依拠する。ノード情報をリンク情報と組み合わせて使用することにより、クラスタリングが収束するまで（クラスタリングアルゴリズムは、レイヤ間で提供される）、クラスタリングされた結果を繰り返し射影し、伝播させる。すなわち、各型の異なる種類のノードおよびリンクが調べられることにより、クラスタリングのために使用することができる構造情報が獲得される。構造情報は、例えば、異なるデータオブジェクトを接続するリンクの型（例えば、リンクが、レイヤ間リンクであるか、またはレイヤ内リンクであるか）を考慮することにより獲得することができる。所定の型のオブジェクトの結果を別の型のオブジェクトのクラスタリング結果に繰り返しクラスタリングすることにより、データの希薄さに関連するクラスタリングの困難を削減することができる。この反復射影により、所定のレイヤクラスタリングにおける類似測度を、別の型のクラスタの個々のグループの代わりのクラスタで計算する。

【 0 0 4 6 】

例えば、2つのオブジェクト型、 $X = \{x_1, x_2, \dots, x_m\}$ および $Y = \{y_1, y_2, \dots, y_n\}$ の観点からプロセスを例示すると、S T S モジュール 2 1 2 はまず、あらゆる従来のクラスタリング方法を使用して、Y 中のオブジェクトを $\{C_1, C_2, \dots, C_k\}$ で表される k 個のクラスタにクラスタリングする。 x X の関係特徴ベクトルを含む M D O ベクトル 2 3 4 は、最初、 $V_x = [v_1, v_2, \dots, v_n]^T$ として定義され、各成分が Y 中の所定のオブジェクトに対応することを想起されたい。Y 中のクラスタにより、各成分が Y の所定のクラスタに対応し、 $x \cdot R_Y C_i$ である場合、 v_i が 0 ではない、 $V_x' = [v_1', v_2', \dots, v_k']^T$ で V_x を置き換える。 v_i' の数値は、オブジェクト x からクラスタ C_i 中の複数のオブジェクトに対する関係の数を表現する、 $|x \cdot R_Y C_i|$ に、または関連するオブジェクトの重要度（オブジェクト重要度は、以下に説明する）に設定することも可能である。したがって、X 中のオブジェクトのクラスタリングは、この新たな型間関係特徴に基づく。プロセスは、所定の型のクラスタリング結果を別の型にそれらの型のレイヤ間関係により繰り返し射影することによって、収束まで続けられる。

【 0 0 4 7 】

以上の強化されたクラスタリングアルゴリズムの利点は、クラスタリング結果が、コンテンツからのデータ分布を反映するだけでなく、他のデータ型との関係も反映することである。このアルゴリズムは、データ希薄性の問題もある程度、解決することができる。固定の特徴空間に関して類似度を定義する既存のクラスタリングアプローチと比較すると、マルチ型データオブジェクトの強化されたクラスタリングを行うための説明するシステムおよび方法は、クラスタリングプロセス中に2つのオブジェクト間で類似度を更新することにより、新たに発見された関係特徴空間に適応する。さらに、所定の実装において、あらゆる従来のクラスタリングアルゴリズムをこの提案するフレームワークに埋め込むことにより、クラスタリングパフォーマンスを向上させることができる。

【0048】

10

(リンク分析およびオブジェクトの重要度)

一部のデータオブジェクトおよびアプリケーションに関して、同一型内の複数のマルチ型データオブジェクト214は、クラスタリングプロセスにおいて異なる重要度を有する可能性がある。通常の実施例には、若干のWebページが権威のあるページであるため、より重要である場合のWebページ/ユーザクラスタリング、および一部のユーザが、アイテムのふさわしさ(belongingness)を決定する際により権威を有すべき場合の、共同フィルタリングなどのためのアイテム/ユーザクラスタリングが含まれる。オブジェクトをノード見なし、オブジェクト間の(複数の)関係をリンクと見なした場合、HITSアルゴリズムなどの従来のリンク分析方法を使用することにより、各データオブジェクトの固有値を計算する。しかし、複数の型のデータオブジェクトが関係する場合、この方法は、異なる型のオブジェクトの重要度が比較できないため、うまく作用しない。

20

【0049】

この問題に対処するため、マルチ型データオブジェクトの強化されたクラスタリングのための説明するシステムおよび方法は、次のとおりHITSアルゴリズムを拡張する。型内でオブジェクト重要度の相互強化を考慮するだけでなく、型間の相互強化も考慮する。各ノードには、ハブスコア(hub score)および権威スコア(authority score)が割り当てられる。

【0050】

簡単にするため、提案するアルゴリズムを例示する実施例として、2つの型の相互に関係のあるオブジェクトを含むケースを使用することを続ける。2つの型のオブジェクト、 $X = \{x_1, x_2, \dots, x_m\}$ 、 $Y = \{y_1, y_2, \dots, y_n\}$ 、ならびに R_x 、 R_y 、 R_{xy} 、および方向性が考慮される場合 R_{yx} という関係を所与として、隣接する行列を使用してリンク情報を表現する。 L_x および L_y がそれぞれ、集合XおよびY内のリンク構造の隣接する行列を表す。 L_{xy} および L_{yx} が、Xの中のオブジェクトからYの中のオブジェクトへのリンクの隣接する行列を表す。例えば、ノード x_i からノード y_j への所定のリンクが存在する場合、 $L_{xy}(i, j) = 1$ である。

30

【0051】

次の2つのレベルの計算が存在する。すなわち、1つのレベルは、同一型からのオブジェクトのハブ値(hub value)および権威値(authority value)が、型内関係によって互いを強化することであり、他方のレベルは、異なる型のノードの重要度が、型間関係によって互いを強化することである。このアプローチにおける計算は、以下のとおり書き表される。

40

【0052】

【数 2】

$$\left\{ \begin{array}{l} a(X) = \beta L_X^T h(X) + (1 - \beta) L_{XY} i(Y) \\ h(X) = \beta L_X a(X) + (1 - \beta) L_{XY} i(Y) \\ i(X) = a(X) + h(X) \\ \\ a(Y) = \gamma L_Y^T h(Y) + (1 - \gamma) L_{YX} i(X) \\ h(Y) = \gamma L_Y a(Y) + (1 - \gamma) L_{YX} i(X) \\ i(Y) = a(Y) + h(Y) \end{array} \right. \quad (2)$$

10

【0053】

ただし、 $a(X)$ および $h(X)$ はそれぞれ、 X 内のノードの権威スコアおよびハブスコアである。同様に、 $a(Y)$ および $h(Y)$ は、 Y 内のノードの権威スコアおよびハブスコアであり、 $i(X)$ および $i(Y)$ はそれぞれ、 X および Y の中のノードの重要度を表す。およびは、異なる関係に由来するリンクの影響を調整する重みパラメータである。

【0054】

20

計算の始めに、すべてのベクトル、 $a(X)$ 、 $h(X)$ 、 $a(Y)$ 、および $h(Y)$ は、1 に初期設定される。ハブスコアおよび権威スコアは、各反復時に等式(2)を使用して更新される。各反復の終了時に、ベクトルは、次の反復計算のために正規化される。このアルゴリズムは、各オブジェクト型内で正規化された一様な重要度を提供し、型間関係を介して他の型の関連するオブジェクトの重要度を考慮することにより、より妥当な結果を獲得する。

【0055】

オブジェクトの重要度スコアを所与として、説明する強化されたクラスタリングプロセスは、オブジェクトの重要度を反映するように変更される。本実装において、 k 平均クラスタリングアルゴリズムが、重み付けされた k 平均アルゴリズムに変更される。すなわち、クラスタ中心を計算する際、クラスタメンバの重み付けされた合計を新たな中心として使用して、クラスタがそれらの重要なオブジェクトの方にバイアスされるようにする。

30

【0056】

以上の観点から、STS モジュール 212 は、複数のマルチ型データオブジェクト 214 の間における型間関係と型内関係の両方に基づき、マルチ型データオブジェクトの重要度を区別する。この重要度が、クラスタリングプロセスに組み込まれる。

【0057】

(入札の語の例示的处理)

エンドユーザ(例えば、広告主、Web サイト主催者など)から(複数の)語 208 を受け取ったことに応答して、STS モジュール 212 は、(複数の)語 208 を複数の強化された語のクラスタ 236 の中の語/句のそれぞれの語/句と比較する。複数の強化された語のクラスタ 236 は、Web ページおよびユーザとの語の相互関係に由来する、文脈上で互いに関係するだけでなく、意味上でも互いに関係する複数の語を含むので、(複数の)語 208 は、複数の関係のある履歴上の文脈、すなわち「意味」の観点から評価される。

40

【0058】

所定の実装において、STS モジュール 212 により、(複数の)語 208 が、強化されたクラスタ 236 からの(複数の)語と合致すると決定する場合、検索語提案モジュール 212 は、強化されたクラスタ 236 から提案される語のリスト 210 を生成する。本実装において、合致は、正確な合致であっても、単数形/複数形、つづり間違い、句読点

50

などの少数の違いを伴う合致であってもよい。戻されるリストは、F O Oと信頼度値の結合により順に並べられる。

【 0 0 5 9 】

所定の実装において、(複数の)語がクラスタからの語に合致する場合、そのクラスタが、提案される語のリストの中でエンドユーザに戻される。提案される語のリスト210は、(複数の)語に意味上、および/または文脈上、関係があると決定された語/句、それぞれの(複数の)語と(複数の)語との類似測度(信頼度値)、およびそれぞれの(複数の)語出現頻度(F O O)を含む。戻されるリスト210は、F O Oと信頼度値の結合により順に並べられる。

【 0 0 6 0 】

S T Sモジュール212により、語208が複数の強化された語のクラスタ236の中の語と合致すると決定する場合、検索語提案モジュール212により、複数の強化された語のクラスタ236の複数のクラスタの中の語から、複数の提案された語リスト210を生成する。複数のリストは、クラスタサイズにより順に並べられ、各リスト内の語は、F O Oと信頼度値の結合により順に並べられる。

【 0 0 6 1 】

合致するクラスタがないと識別された場合、クエリの語は、低いF O Oを有するクエリの語から生成された、拡張されたクラスタ群に対してさらにマッチングされる。所定の実装において、低いF O Oを有するクエリ語は、高い出現頻度の履歴クエリログの語から生成された、複数の強化された語のクラスタ236に関する分類子(classifier)(例えば、K-最近傍(nearest neighbor)の分類子)を訓練することによってクラスタリングされる。低い出現頻度を有すると決定された履歴クエリの語が1つずつ、検索エンジンに送信される。次に、戻された検索結果の選定された検索結果(例えば、第1の最上位ランクのWebページ、および/または等)から、特徴が抽出される。抽出された特徴を、正規化し、低いF O Oを有するクエリの語を表現するのに使用する。次に、クエリの語は、既存の複数のクラスタに分類することにより、訓練された分類子に基づいて複数の拡張されたクラスタが生成される。次に、エンドユーザが送信した(複数の)語を、これらの拡張されたクラスタの観点から評価することにより、提案される語のリストを識別し、エンドユーザに戻す。

【 0 0 6 2 】

(低いF O O語の分類)

高いF O O(出現頻度)のクエリの語222から生成され、強化された語の複数のクラスタ236が、エンドユーザが入力した(複数の)語208と同一の語を含まない場合、S T Sモジュール212により、高いF O O(出現頻度)のクエリログの語222から生成され、強化された語の複数のクラスタ236から、訓練された分類子238を生成する。強化された語の複数のクラスタ236の中の語は、分類操作に適切なベクトル空間モデルの中で対応するキーワードベクトルを既に有する。さらに、ストップワードの削除および単語のステミング(stemming)(接尾辞削除)により、(複数のクラスタ236に基づく)語のベクトル234の次元が縮小される。所定の実装において、さらなる次元縮小技術、例えば、特徴選択またはパラメータ変更(re-parameterization)を使用することができる。

【 0 0 6 3 】

本実装において、クラスが未知のクエリの語222を分類するため、S T Sモジュール212は、k-最近傍の分類子(k-Nearest Neighbor classifier)のアルゴリズムを使用することにより、クラスが既知のすべてのクエリの語222に対応する特徴ベクトルに依拠する、クラスが既知のすべてのクエリの語222の中でk個の最も類似した近傍を求め、近傍のクラスラベルの重み付けされたマジョリティ(weighted majority)を使用することにより、新たなクエリの語のクラスを予測する。この場合、既に強化された語の複数のクラスタ236の中にあるおののクエリの語には、対応するクラスタのラベルと同一のラベルが割り当てられる一方で、

10

20

30

40

50

おのこの強化されたクラスタ 2 3 6 には、単なるシーケンス番号でラベルが付けられる。それらの近傍には、X に対する各近傍の類似度を使用して重みが付けられ、ただし、類似度は、2 つのベクトルの間のユークリッド距離、またはコサイン値で測定される。コサイン類似度は、以下のとおりである。すなわち、

【 0 0 6 4 】

【 数 3 】

$$\text{sim}(X, D_j) = \frac{\sum_{t_i \in (x \cap D_j)} x_i \cdot d_{ij}}{\|X\|_2 \cdot \|D_j\|_2}$$

10

【 0 0 6 5 】

ただし、X は、試験語 (t e s t t e r m)、すなわち、ベクトルとして表現される、分類されるべきクエリの語であり、D_j は、第 j 番の訓練中の語であり、t_i は、X と D_j が共有する語であり、x_i は、X 中のキーワード t_i の重みであり、d_{ij} は、D_j 中のキーワード t_i の重みであり、

【 0 0 6 6 】

【 数 4 】

$$\|X\|_2 = \sqrt{x_1^2 + x_2^2 + x_3^2}$$

20

【 0 0 6 7 】

は、X のノルムであり、 $\|D_j\|_2$ は、D_j のノルムである。したがって、試験語 X のクラスラベルは、以下のとおり、すべての近傍のクラスラベルの重み付けされたマジョリティである。すなわち、

【 0 0 6 8 】

【 数 5 】

$$\text{label}(X) = \arg \max_{l_i} \left(\sum_{\text{ラベル}(D_j)=l_i \text{ である場合のすべての } D_j} \text{sim}(X, D_j) \right)$$

【 0 0 6 9 】

30

別の実装において、最近傍の分類技術以外の異なる統計分類 - 機械学習技術 (例えば、回帰モデル、ベイズ分類子、判定ツリー、ニューラルネットワーク、およびサポートベクトルマシンを含む) を使用して、訓練された分類子 2 3 8 が生成される。

【 0 0 7 0 】

S T S モジュール 2 1 2 により、低い F O O (出現頻度) のクエリの語 2 2 4 を 1 つずつ (それぞれ検索クエリ 2 2 7 を介して)、検索エンジン 2 2 8 に送信する。特定の検索クエリ 2 2 7 に関連する (複数の) 検索結果 2 3 0 を受け取ったことに応答し、既に説明した技術を使用して、S T S モジュール 2 1 2 は、(複数の) 検索結果 2 3 0 によって識別された 1 つまたは検索された検索結果 2 3 0 から、特徴を抽出する (抽出された特徴 2 3 2)。本実装において、特徴は、第 1 の最上位ランクの (複数の) 検索結果 2 3 0 から抽出される。おのこの検索され、解析された (複数の) 検索結果 2 3 0 に関して、S T S モジュール 2 1 2 は、抽出された特徴 2 3 2 のそれぞれのレコードの中に以下の情報を格納する。すなわち、U R L、結果タイトル、クエリの語の短い説明および / または文脈、ならびに検索された検索結果 2 3 0 を獲得するのに使用された検索クエリ 2 2 7 である。次に、S T S モジュール 2 1 2 は、低い F O O クエリの語 2 2 4 に由来する抽出済みの特徴 2 3 2 をトークン化し、次元を縮小し、正規化することにより、語の複数のベクトル 2 3 4 を生成する。次に、S T S モジュール 2 1 2 は、クエリの語を、語のクラスタ 2 3 6 のそれぞれの集合にクラスタリングする。このクラスタリング操作は、(高い F O O クエリの語 2 2 2 から生成された) 訓練された分類子 2 3 8 を使用して実行される。

40

【 0 0 7 1 】

50

S T S モジュール 2 1 2 は、(低い F O O クエリの語 2 2 4 に基づいて生成された、) これらの拡張された語の複数のクラスタの観点からエンドユーザが送信した(複数の) 語 2 0 8 を評価することにより、1 つまたは複数の提案される語のリスト 2 1 0 を識別し、エンドユーザに戻す。例示的なそのような手続きを段落 (0 0 6 3) および段落 (0 0 6 6) で上述しており、以下のセクションで説明する。

【 0 0 7 2 】

(例示的な検索語提案のリスト)

提案される語のリスト 2 1 0 は、例えば、(複数の) 語 2 0 8 に関係があると決定された語、それぞれの(複数の) 語と(複数の) 語 2 0 8 との類似測度(信頼度値)、およびそれぞれの(複数の) 語の F O O (出現頻度)、すなわち、履歴クエリログ内の頻度を含む。関係のある(複数の) 語を識別する、類似測度を生成する、および F O O 値を生成するための技術は、上述した。

10

【 0 0 7 3 】

表 1 は、「 m a i l 」という語 2 0 8 に関係があると決定された語の例示的な提案される語のリスト 2 1 0 を示す。本実施例では、(複数の) 語 2 0 8 に関係のある語が、「提案される語」という題名が付けられた第 1 列の中に示されている。

【 0 0 7 4 】

【表 1】

表 1

入札の語「MAIL」に関する例示的な提案される語のリスト

提案される語	類似度	頻度	<文脈>
hotmail	0.246142	93161	関係のあるオンライン電子メール
yahoo	0.0719463	165722	
mail.com	0.352664	1455	
yahoo mail	0.0720606	39376	
www.mail.com	0.35367	711	
email.com	0.484197	225	
www.hot	0.186565	1579	
www.msn.com	0.189117	1069	
mail.yahoo.com	0.0962268	4481	
free email	0.230611	1189	
www.aolmail.com	0.150844	654	
check mail	0.221989	66	
check email	0.184565	59	
msn passport	0.12222	55	
www.webmail.aol.com	0.0200538	108	
webmail.yahoo.com	0.08789	71	
free email account	0.0234481	65	
提案される語	類似度	頻度	
mail	1	2191	関係のある従来のメール
usps	0.205141	4316	
usps.com	0.173754	779	
united parcel service	0.120837	941	
postal rates	0.250423	76	
stamps	0.156702	202	
stamp collecting	0.143618	152	
state abbreviations	0.104614	300	
postal	0.185255	66	
postage	0.180112	55	
postage rates	0.172722	51	
usps zip codes	0.138821	78	
us postmaster	0.109844	58	

【0075】

表 1 を参照して、提案される語のリストの中の語が、語類似度値（「類似度」という題名が付いた第 2 列を参照）および出現頻度スコア（「頻度」という題名が付いた第 3 列を参照）にマッピングされていることに留意されたい。「語のクラスタリング」という題名が付いたセクションにおいて後述するとおり計算される各語類似度値は、対応する提案される語（第 1 列）と、本実施例では「mail」である（複数の）語 208 との間の類似測度を提供する。各頻度値、または頻度スコアは、履歴クエリログ内で提案される語が出

10

20

30

40

50

を行うための図3における例示的处理手順300の続きである。説明の目的のために、処理手順の操作は、図2の特徴に関連して説明する。ブロック402では、エンドユーザから(複数の)語208(図2)を受け取ったことに応答して、STSモジュール212により、(複数の)語208に大いに類似し、関係があると決定された、強化された語の複数のクラスタ236からのあらゆる(複数の)語から、提案される語のリスト210を生成する。異なるオブジェクト型の間の相互関係を使用することにより、クラスタリングを向上させる。ブロック404では、STSモジュール212により、キーワードクラスタ236からのあらゆる(複数の)語が、(複数の)語208に大いに類似する/関係のがあると決定されたかどうかを決定する。類似する/関係のがあると決定された場合、処理手順は、ブロック406に続き、STSモジュール212により、対応する提案される語のリスト210をエンドユーザに送信する。類似しない/関係がない場合、処理手順は、ページ上の参照符号「B」で示すとおり、図5のブロック502に継続する。

10

【0081】

図5は、検索語提案のためにマルチ型データオブジェクトの強化されたクラスタリングを行うための、図3および図4における例示的处理手順300の続きである。説明の目的のために、処理手順の操作は、図2の特徴に関連して説明する。(すべての参照符号は、構成要素が最初に導入される図面の番号で始める)。ブロック502では、STSモジュール212により、強化された語の複数のクラスタ236から分類子238(訓練された分類子)を生成し、複数のクラスタ236は、この時点で、高い出現頻度のクエリの語222に基づいている。ブロック504では、STSモジュール212により、低い出現頻度のクエリの語224を1つずつ検索エンジン228に送信し、対応する検索結果230を受け取る。ブロック506では、STSモジュール212により、検索結果230からスニペット記述(抽出された特徴232)を抽出し、その記述から語の複数のベクトル234を生成する。

20

【0082】

ブロック508では、STSモジュール212により、訓練された分類子238の観点から、低い出現頻度のクエリの語224から生成された語の複数のベクトル234を分類することにより、低い出現頻度のクエリの語224に基づくそれぞれの強化された語の複数のクラスタ236を生成する。ブロック510では、STSモジュール212により、(複数の)語208と大いに類似すると決定された低い出現頻度のクエリの語224に基づく強化された語の複数のクラスタ236からのキーワード/重要句から、提案される語のリスト210を生成する。ブロック512では、STSモジュール212により、提案される語のリスト210をエンドユーザに送信する。

30

【0083】

図6は、図3のブロック312の強化されたクラスタリング操作の例示的な詳細を示す。説明の目的のために、ブロック310の操作は、図1および図2の特徴に関連して説明する。STSモジュール212により実装される強化されたクラスタリングアルゴリズムへの入力には、グラフ100のようなマルチレイヤのフレームワークグラフが含まれ、グラフ100は、識別され、重みが付けられたオブジェクト間およびオブジェクト内関係を含む、対応する複数のノードのコンテンツ特徴 f_i および g_j を含む。このクラスタリングアルゴリズムの出力には、マルチ型データオブジェクトの強化されたクラスタリングを反映する新たなフレームワークグラフ100が含まれる。新たなフレームワークグラフの若干の実装において、新たなノード位置に変更され、および/またはグラフ100の次元を縮小するために別のノードとマージされた各々の古いノードの変化を例示することが可能である。

40

【0084】

ブロック602では、元のフレームワークグラフが(各クラスタリング反復の前に)入力される。ブロック604では、考慮されている各ノードの重要度が、等式(2)を使用して決定または計算される。ブロック606では、任意のレイヤが、クラスタリングのために選択される。ブロック608では、選択されたレイヤの中のノードを適切となるよう

50

にクラスタリングすることにより（例えば、コンテンツ特徴に従って）、強化されたクラスタ 236 が生成される。若干の実装において、ノードは、所望のフィルタリングアルゴリズム（図示せず）を使用してフィルタ処理することにより、クラスタリングを向上させることができる。ブロック 610 では、各クラスタの複数のノードが 1 つのノードにマージされる。例えば、フィルタ処理の後に 2 つの候補ノードが存在する場合、最も近接した 2 つの候補ノードを、例えば、その 2 つの候補ノードのベクトル値を平均することにより、マージすることができる。このマージにより、個々のノードを結合することができることによって、考慮されなければならないノードの数を減らすことが可能になる。真の意味において、マージする操作を使用することにより、重複または重複に近いノード（near-duplicate）の出現を減らすことができる。ブロック 612 では、610 におけるマージに基づいて、対応するリンクが更新される。ブロック 614 では、クラスタリングアルゴリズムは、クラスタリングのために第 2 のレイヤに（任意に選択されたレイヤから）切り替える。ブロック 312 の操作は、ページ上の参照符号「C」で示されるとおり図 7 のブロック 702 に継続する。

【0085】

図 6 の操作を参照すると、最初のクラスタリングパス（pass）において、コンテンツ特徴だけが利用されることに留意されたい。ほとんどのケースでは、リンク特徴は、当初は、希薄すぎてクラスタリングに役立たない。図 7 を参照して以下に説明する、後続のクラスタリングパスでは、コンテンツ特徴とリンク特徴とを結合することにより、クラスタリングの有効性が高められる。コンテンツ特徴とリンク特徴とを結合することにより、重みを異なる値で指定し、結果を比較することが可能であり、向上した精度を有するクラスタリングを提供することが可能である。

【0086】

図 7 は、図 3 および図 6 のブロック 312 の強化されたクラスタリング操作の例示的な続きを示す。ブロック 702 では、第 2 のレイヤの複数のノードが、複数のノードのコンテンツ特徴および更新されたリンク特徴に従ってクラスタリングされる。ブロック 704 では、各クラスタの複数のノードが 1 つのノードにマージされる。ブロック 706 では、他方のレイヤの元のリンク構造および元の複数のノードが復元される。ブロック 708 では、第 2 のレイヤの各クラスタの複数のノードがマージされ、対応するリンクが更新される。ブロック 710 では、この反復クラスタリングプロセスが、コンピュータ環境内で継続される。ブロック 712 では、フレームワークグラフ 100 の改訂されたバージョンを出力する。

【0087】

（例示的動作環境）

図 8 は、検索語提案のためにマルチ型データオブジェクトの強化されたクラスタリングを行うための、図 2 におけるシステム 200、および図 3 から図 6 の方法を完全に、または部分的に実装することができる適切なコンピュータ処理環境 800 の例を例示する。例示的コンピュータ処理環境 800 は、適切なコンピュータ処理環境の一例に過ぎず、本明細書で説明するシステムおよび方法の用法または機能の範囲について何ら限定を示唆することを意図するものではない。また、コンピュータ処理環境 800 は、コンピュータ処理環境 800 に例示したコンポーネントのいずれの 1 つ、または組み合わせに関連する依存関係または要件も有するものと解釈してはならない。

【0088】

本明細書で説明する方法およびシステムは、他の多数の汎用または専用の、コンピュータ処理システムの環境または構成で機能する。使用に適する可能性がある周知のコンピュータ処理のシステム、環境、および/または構成の例には、パーソナルコンピュータ、サーバコンピュータ、マルチプロセッサシステム、マイクロプロセッサベースのシステム、ネットワーク PC、ミニコンピュータ、メインフレームコンピュータ、以上のシステムまたは装置のいずれかを含む分散コンピュータ処理環境などが含まれるが、以上には限定されない。また、フレームワークのコンパクトまたはサブセットのバージョンを、ハンドヘ

ルドコンピュータ、または他のコンピュータ処理装置などの、限られたリソースのクライアントにおいて実装することもできる。本発明は、通信ネットワークを介してリンクされたりリモート処理ユニット群によってタスクが実行される分散コンピュータ処理環境において実施される。分散コンピュータ処理環境では、プログラムモジュール群は、ローカルおよびリモートのメモリ記憶装置の両方の中に配置することができる。

【0089】

図8を参照すると、検索語提案のためにマルチ型データオブジェクトの強化されたクラスタリングを行うための例示的システムは、コンピュータ810の形態で汎用コンピュータ処理装置を含む。コンピュータ810の以下に説明する態様は、クライアントコンピュータ処理装置PSSサーバ202(図2)および/またはクライアントコンピュータ処理装置206の例示の実装である。コンピュータ810のコンポーネントには、処理ユニット(群)820、システムメモリ830、ならびにシステムメモリから処理ユニット820までを含む様々なシステムコンポーネントを結合するシステムバス821を含むことができるが、以上には限定されない。システムバス821は、メモリバスまたはメモリコントローラ、周辺バス、およびあらゆる様々なバスアーキテクチャを使用するローカルバスを含め、いくつかの型のバス構造のいずれであることも可能である。例として、限定としてではなく、そのようなアーキテクチャには、インダストリスタンダードアーキテクチャ(ISA)バス、マイクロチャネルアーキテクチャ(MCA)バス、エンハンストISA(EISA)バス、ビデオエレクトロニクススタンダーズアソシエーション(VESA)ローカルバス、およびメザニン(Mezzanine)バスとしても知られるペリフェラルコンポーネントインターコネクト(PCI)バスを含むことが可能である。

【0090】

コンピュータ810は、通常、様々なコンピュータ読取可能な媒体を含む。コンピュータ読取可能な媒体は、コンピュータ810によりアクセスすることができるあらゆる利用可能な媒体であることが可能であり、揮発性および不揮発性の媒体、取り外し可能および固定の媒体がともに含まれる。例として、限定としてではなく、コンピュータ読取可能な媒体は、コンピュータ記憶媒体および通信媒体を備えることが可能である。コンピュータ記憶媒体には、コンピュータ読取可能な命令、データ構造、プログラムモジュール、または他のデータなどの情報を格納するためにあらゆる方法または技術で実装された、揮発性および不揮発性の、取り外し可能および固定の媒体が含まれる。コンピュータ記憶媒体には、RAM、ROM、EEPROM、フラッシュメモリもしくは他のメモリ技術、CD-ROM、デジタル多用途ディスク(DVD)もしくは他の光ディスクストレージ、磁気カセット、磁気テープ、磁気ディスクストレージもしくは他の磁気記憶装置、または所望の情報を格納するのに使用することができおよびコンピュータ810によりアクセスすることができるあらゆる他の媒体が含まれるが、以上には限定されない。

【0091】

通信媒体は、通常、コンピュータ読取可能な命令、データ構造、プログラムモジュールまたは搬送波もしくは他の移送機構などの変調されたデータ信号の中の他のデータを含み、およびあらゆる情報配信媒体が含まれる。「変調されたデータ信号」という用語は、信号内に情報を符号化するような形で、特性の1つまたは複数が設定または変更された信号を意味する。例として、限定としてではなく、通信媒体には、有線ネットワークまたは直接有線接続などの有線媒体、ならびに音響、RF、赤外線、および他の無線媒体などの無線媒体が含まれる。以上の媒体のいずれか媒体の組み合わせも、コンピュータ読取可能な媒体の範囲内に含まれなければならない。

【0092】

システムメモリ830は、読み出し専用メモリ(ROM)831およびランダムアクセスメモリ(RAM)832などの、揮発性および/または不揮発性メモリの形態でコンピュータ記憶媒体を含む。始動中などにコンピュータ810内部の要素間で情報を転送するのを助ける基本ルーチンを含む基本入出力システム(BIOS)833が、通常、ROM831の中に格納される。RAM832は、通常、処理ユニット820により即時にアク

セス可能および／または現在処理中のデータおよび／またはプログラムモジュール群を含む。例として、限定としてではなく、図 8 は、オペレーティングシステム 8 3 4、アプリケーションプログラム群 8 3 5、他のプログラムモジュール群 8 3 6、およびプログラムデータ 8 3 8 を例示する。所定の実装において、コンピュータ 8 1 0 は、P S S サーバ 2 0 2 である。本シナリオにおいて、アプリケーションプログラム群 8 3 5 は、検索語提案モデル 2 1 2 を含む。この同一のシナリオでは、プログラムデータ 8 3 8 は、マルチ型データオブジェクト 2 1 4、検索結果 2 3 0、抽出された特徴 2 3 2、M D O ベクトル群 2 3 4、強化されたクラスタ群 2 3 6、訓練された分類子 2 3 8、および他のデータ 2 2 6 を含む。

【 0 0 9 3 】

コンピュータ 8 1 0 は、他の取り外し可能／固定の、揮発性／不揮発性のコンピュータ記憶媒体も含むことが可能である。単に例として、図 8 は、固定の不揮発性の磁気媒体に対して読み出しまたは書き込みを行うハードディスクドライブ 8 4 1、取り外し可能な不揮発性の磁気ディスク 8 5 2 に対して読み出しまたは書き込みを行う磁気ディスクドライブ 8 5 1、および C D - R O M または他の光媒体などの取り外し可能な不揮発性の光ディスク 8 5 6 に対して読み取りまたは書き込みを行う光ディスクドライブ 8 5 5 を例示する。例示的動作環境において使用することができる他の取り外し可能な／固定の、揮発性／不揮発性のコンピュータ記憶媒体には、磁気テープカセット、フラッシュメモリカード、デジタル多用途ディスク、デジタルビデオテープ、固体素子 R A M、固体素子 R O M などが含まれるが、以上には限定されない。ハードディスクドライブ 8 4 1 は、通常、インタフェース 8 4 0 などの固定のメモリインタフェースを介してシステムバス 8 2 1 に接続され、ならびに磁気ディスクドライブ 8 5 1 および光ディスクドライブ 8 5 5 は、通常、インタフェース 8 5 0 などの取り外し可能なメモリインタフェースによりシステムバス 8 2 1 に接続される。

【 0 0 9 4 】

以上に説明し、図 8 に例示した駆動装置、および関連するコンピュータ記憶媒体により、コンピュータ読取可能な命令、データ構造、プログラムモジュール、および他のデータのストレージがコンピュータ 8 1 0 に提供される。図 8 では、例えば、ハードディスクドライブ 8 4 1 が、オペレーティングシステム 8 4 4、アプリケーションプログラム群 8 4 5、他のプログラムモジュール群 8 4 6、およびプログラムデータ 8 4 8 を格納しているのを例示する。これらのコンポーネントは、オペレーティングシステム 8 3 4、アプリケーションプログラム群 8 3 5、他のプログラムモジュール群 8 3 6、およびプログラムデータ 8 3 8 と同一であることも、異なることも可能であることに留意されたい。オペレーティングシステム 8 4 4、アプリケーションプログラム群 8 4 5、他のプログラムモジュール群 8 4 6、およびプログラムデータ 8 4 8 に、本明細書では、それらが少なくとも異なるコピーであることを例示するために異なる参照符号を付ける。

【 0 0 9 5 】

ユーザは、キーボード 8 6 2 および、マウス、トラックボール、またはタッチパッドと一般に呼ばれるポインティング装置 8 6 1 などの入力装置群を介して、コマンドおよび情報をコンピュータ 8 1 0 に入力することができる。他の入力装置群（図示せず）には、マイクロフォン、ジョイスティック、ゲームパッド、衛星受信アンテナ、スキャナなどを含むことが可能である。以上および他の入力装置群は、システムバス 8 2 1 に結合されたユーザ入力インタフェース 8 6 0 を介して処理ユニット 8 2 0 に接続される場合が多いが、パラレルポート、ゲームポート、またはユニバーサルシリアルバス（U S B）などの、他のインタフェースおよびバス構造により接続してもよい。

【 0 0 9 6 】

また、モニタ 8 9 1 または他の型のディスプレイ装置も、ビデオインタフェース 8 9 0 のようなインタフェースを介して、システムバス 8 2 1 に接続される。モニタに加えて、コンピュータは、出力周辺インタフェース 8 9 5 を介して接続することができるスピーカ 8 9 8 やプリンタ 8 9 6 などの、他の周辺出力装置群も含むことが可能である。

【 0 0 9 7 】

コンピュータ 8 1 0 は、リモートコンピュータ 8 8 0 など、1 つまたは複数のリモートコンピュータへの論理接続を使用するネットワーク化された環境において動作する。リモートコンピュータ 8 8 0 は、パーソナルコンピュータ、サーバ、ルータ、ネットワーク PC、ピア装置、または他の共通のネットワークノードであることが可能であり、コンピュータ 8 8 0 の特定の実装との相関的要素として、コンピュータ 8 1 0 に関係のある上述した要素の多くまたはすべてを含むことが可能であるが、メモリ記憶装置 8 8 1 だけを図 8 に例示している。図 8 に示す論理接続には、ローカルエリアネットワーク (LAN) 8 8 1 およびワイドエリアネットワーク (WAN) 8 8 3 を含むが、他のネットワークも含むことが可能である。そのようなネットワーキング環境は、オフィス、企業全体のコンピュータネットワーク、イントラネット、およびインターネットで一般的に見られる。

10

【 0 0 9 8 】

LAN ネットワーキング環境において使用される場合、コンピュータ 8 1 0 は、ネットワークインタフェースまたはネットワークアダプタ 8 8 0 を介して LAN 8 8 1 に接続される。WAN ネットワーキング環境において使用される場合、コンピュータ 8 1 0 は、通常、モデム 8 8 2、またはインターネットなどの WAN 8 8 3 を介して通信を確立するための他の手段を含む。モデム 8 8 2 は、内蔵型でも外付け型でもよく、ユーザ入力インタフェース 8 6 0、または他の適切な機構を介してシステムバス 8 2 1 に接続することができる。ネットワーク化された環境では、コンピュータ 8 1 0 との関係を示したプログラムモジュール群、またはプログラムモジュール群の一部は、リモートメモリ記憶装置の中に格納することができる。例として、限定としてではなく、図 8 は、リモートアプリケーションプログラム群 8 8 5 がメモリ装置 8 8 1 上に存在するのを例示する。示したネットワーク接続は例示的であり、コンピュータ間で通信リンクを確立する他の手段も使用することができる。

20

【 0 0 9 9 】

(結 論)

検索語提案のためにマルチ型データオブジェクトの強化されたクラスタリングを行うためのシステムおよび方法を、構造上の特徴、および / または方法上の操作もしくは動作に特有の言い回しで説明してきたが、添付の特許請求の範囲において定義する実装は、説明した特定の特徴または動作に必ずしも限定されないことを理解されたい。例えば、マルチ型データオブジェクトの強化されたクラスタリングを、検索語提案のアプリケーションに関して説明したが、マルチ型データオブジェクトの強化されたクラスタリングは、クラスタリングを利用する他の多くの型のアプリケーションにも適用することができる。したがって、特定の特徴および動作を、請求の対象を実装するための例示的形態として開示する。

30

【 図面の簡単な説明 】

【 0 1 0 0 】

【 図 1 】 異質なデータオブジェクト / ノードの複数のレイヤ 1 0 2、ならびに関連するレイヤ間およびレイヤ内データオブジェクトリンク / 関係を含むマルチレイヤフレームワークグラフ 1 0 0 を示す図である。

40

【 図 2 】 検索語提案のためにマルチ型データオブジェクトの強化されたクラスタリングを行うための例示的システムを示す図である。

【 図 3 】 検索語提案のためにマルチ型データオブジェクトの強化されたクラスタリングを行うための例示的处理手順を示す図である。

【 図 4 】 検索語提案のためにマルチ型データオブジェクトの強化されたクラスタリングを行うための図 3 の例示的处理手順 3 0 0 の続きを示す図である。

【 図 5 】 検索語提案のためにマルチ型データオブジェクトの強化されたクラスタリングを行うための図 3 および図 4 の例示的处理手順 3 0 0 の続きを示す図である。

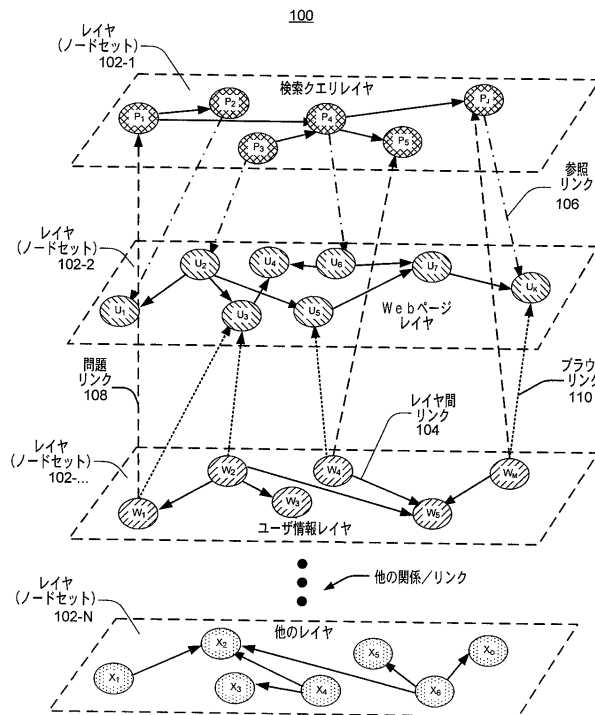
【 図 6 】 図 3 のブロック 3 1 2 の強化されたクラスタリング操作の例示的詳細を示す図である。

50

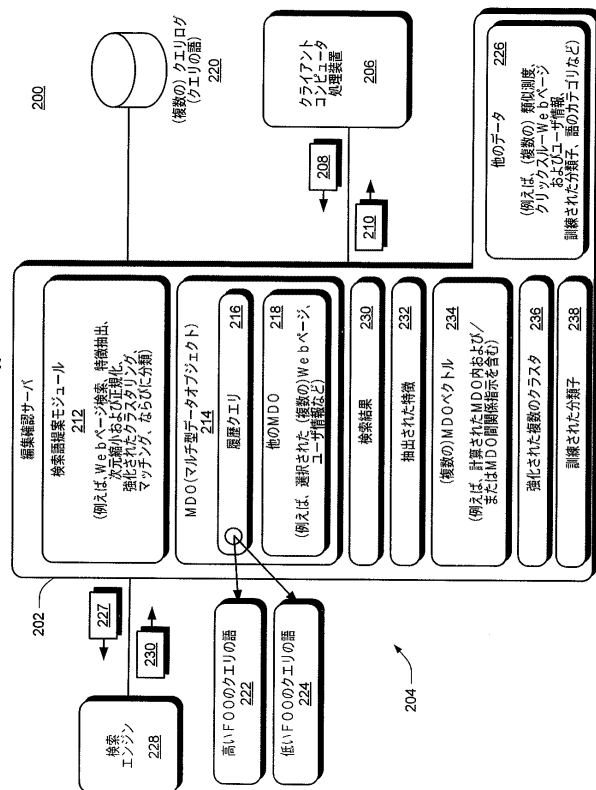
【図 7】図 3 のブロック 3 1 2 および図 6 の強化されたクラスタリング操作の例示的続きを示す図である。

【図 8】検索語提案のためにマルチ型データオブジェクトの強化されたクラスタリングを行うための上述のシステム、装置、および方法を完全に、または部分的に実装することができる適切な例示的コンピュータ処理環境を示す図である。

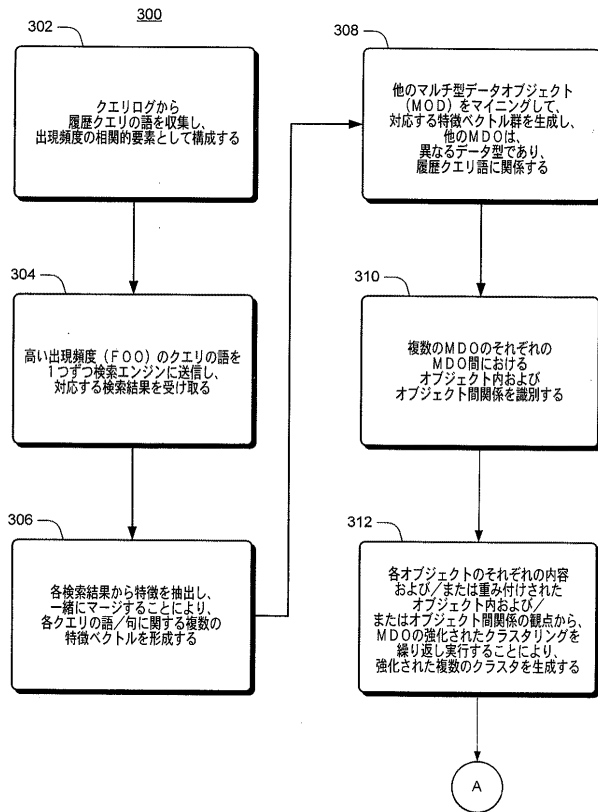
【図 1】



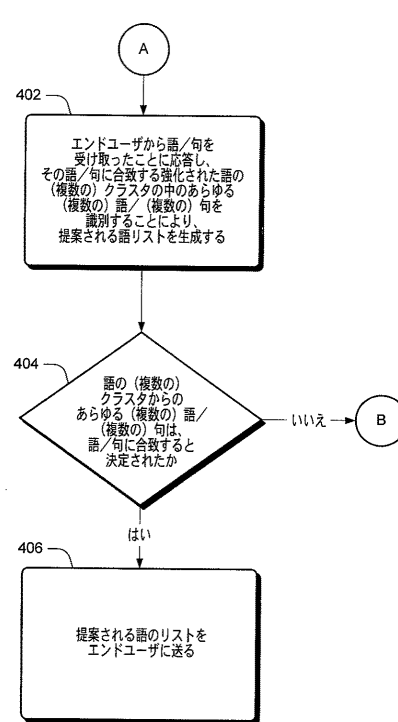
【図 2】



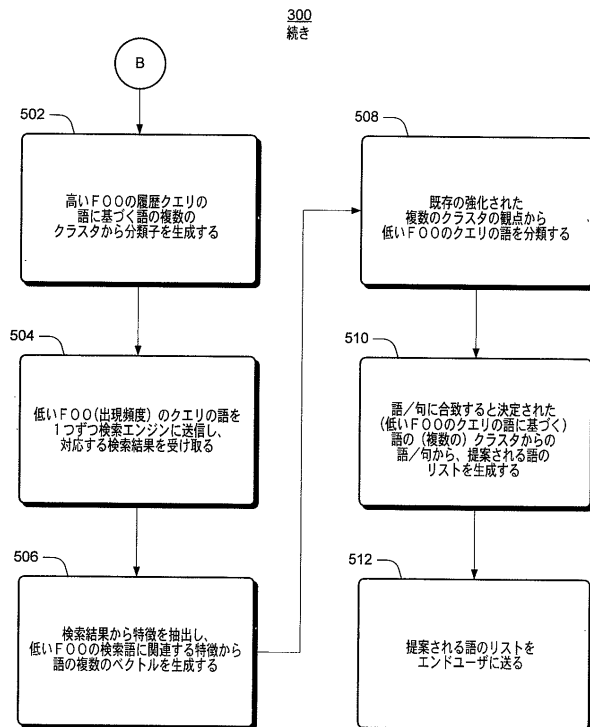
【図 3】



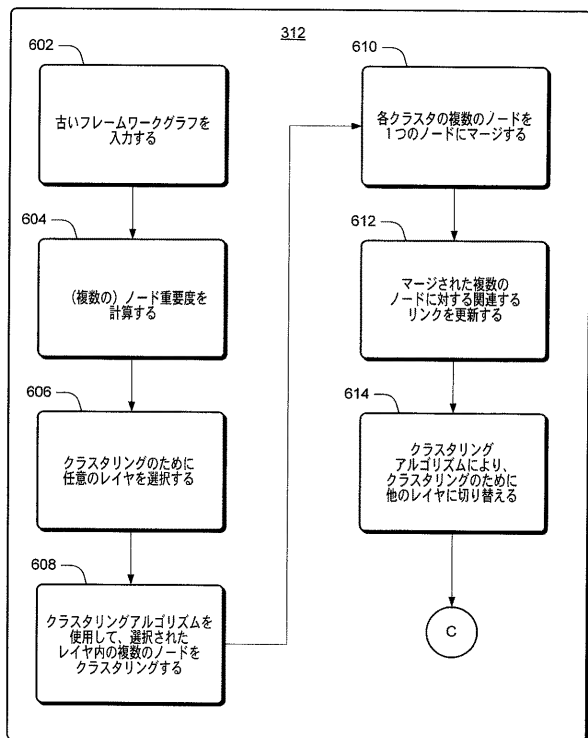
【図 4】



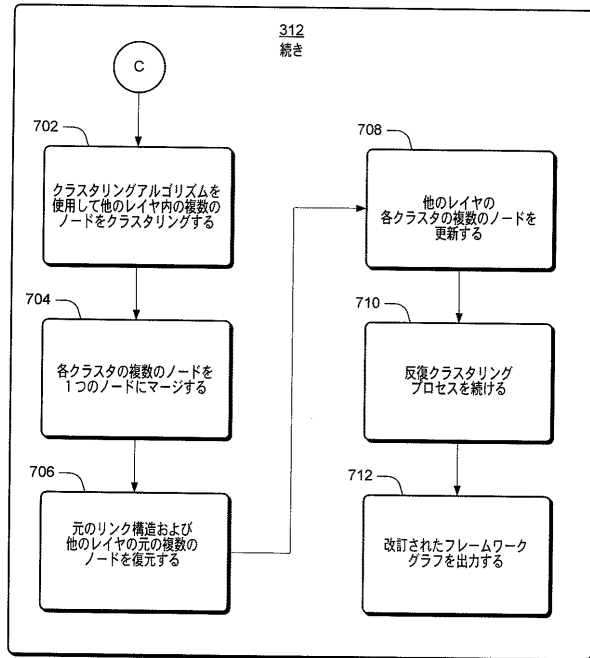
【図 5】



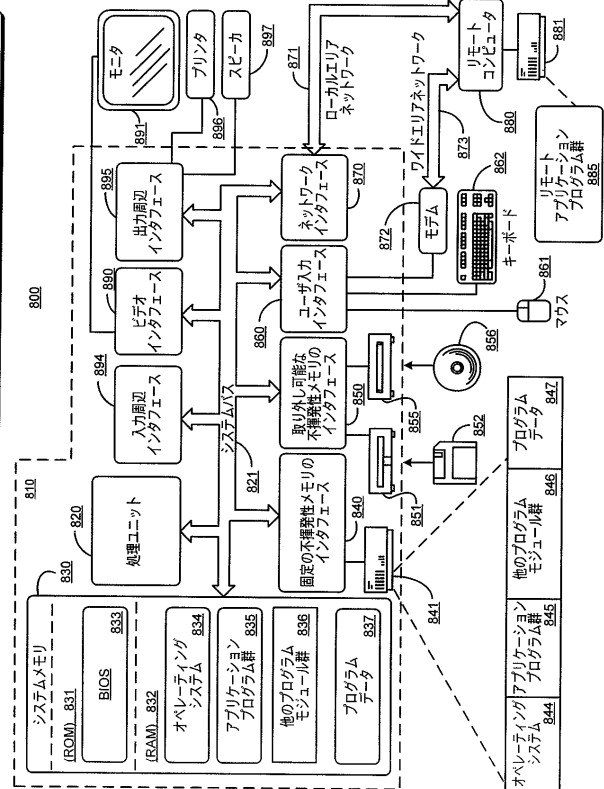
【図 6】



【図 7】



【図 8】



フロントページの続き

- (72)発明者 ホア - ジュン チェン
アメリカ合衆国 98052 ワシントン州 レッドモンド ワン マイクロソフト ウェイ マ
イクロソフト コーポレーション内
- (72)発明者 リー リー
アメリカ合衆国 98052 ワシントン州 レッドモンド ワン マイクロソフト ウェイ マ
イクロソフト コーポレーション内
- (72)発明者 タレック ナジム
アメリカ合衆国 98052 ワシントン州 レッドモンド ワン マイクロソフト ウェイ マ
イクロソフト コーポレーション内
- (72)発明者 ウェイ - イェン マ
アメリカ合衆国 98052 ワシントン州 レッドモンド ワン マイクロソフト ウェイ マ
イクロソフト コーポレーション内
- (72)発明者 イェン リー
アメリカ合衆国 98052 ワシントン州 レッドモンド ワン マイクロソフト ウェイ マ
イクロソフト コーポレーション内
- (72)発明者 チェン ツェン
アメリカ合衆国 98052 ワシントン州 レッドモンド ワン マイクロソフト ウェイ マ
イクロソフト コーポレーション内

審査官 岩間 直純

- (56)参考文献 特開2002-215674(JP, A)
BEEFERMAN D, BERGER A, Agglomerative clustering of a search engine query log, Proceedi
ngs of the sixth ACM SIGKDD international conference on Knowledge discovery and data m
ining, 米国, ACM, 2000年, pp.407-416
WEN J.R, et al., Query clustering using user logs, ACM Transactions on Information Sys
tems (TOIS), 米国, ACM, 2002年 1月, Vol.20, No.1, pp.59-81

- (58)調査した分野(Int.Cl., DB名)
G06F 17/30