

(19) 日本国特許庁(JP)

(12) 特許公報(B2)

(11) 特許番号

特許第4726528号
(P4726528)

(45) 発行日 平成23年7月20日 (2011.7.20)

(24) 登録日 平成23年4月22日 (2011.4.22)

(51) Int. Cl. F I
G06F 17/30 (2006.01)
 G06F 17/30 320D
 G06F 17/30 210D
 G06F 17/30 170A

請求項の数 43 外国語出願 (全 26 頁)

(21) 出願番号	特願2005-118051 (P2005-118051)	(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2005-302042 (P2005-302042A)		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11日 (2008.4.11)		クロソフト ウェイ
(31) 優先権主張番号	10/825,894	(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】 マルチセンスクエリに関する関連語提案

(57) 【特許請求の範囲】

【請求項1】

関連語提案のためにプロセッサにより実行されるコンピュータ実行可能方法であって、
マルチセンスクエリを介して、検索結果を調べるステップであって、
前記マルチセンスクエリが、以下のステップ、すなわち、
サブミットした用語／フレーズに意味的に関連している用語／フレーズを判定するステ
ップであって、意味関係が、意味を判定するために、用語／フレーズの内容を検索するこ
とにより、発見されるステップと、
出現度数（F O O）のしきい値を構成するステップと、
前記構成されたしきい値に基づき、高F O O、または、低F O Oを履歴クエリに割り振
るステップと、
検索エンジンに以前にサブミットされた1組の高出現度数（F O O）履歴クエリに関連
する検索結果から用語ベクトルを生成するステップと、
用語ベクトルの計算された類似性に応じて用語クラスタを生成するステップであって、
前記計算された類似性である $sim(q_j, q_k)$ は、以下の式、すなわち、

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

に従って決定されるステップであって、

式中で、 d は、ベクトル次数を表し、 q は、クエリを表し、 k は、次数インデックスであり、

i 番目のベクトルの j 番目の用語についての重み w は、以下の式、すなわち

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

に従って計算され、

式中で、 TF_{ij} は、用語頻度を表し、 N は、クエリ用語の合計数であり、 DF_j は、 i 番目のベクトルの j 番目の用語を含む抽出済みのフィーチャレコード数である前記決定されるステップと

を含む前記検索結果を調べるステップと、

エンティティから用語／フレーズを受信したことに応答して、1つまたは複数の関連語提案を識別するために、前記用語クラスタ中の用語／フレーズを考慮して、前記マルチセンスクエリを介して前記用語／フレーズを評価するステップであって、前記識別することは、 FOO と確信度値との組み合わせに基づいて行われることを特徴とする評価するステップと、

前記 FOO と確信度値との組み合わせにより順序付けられる少なくとも1つの提案された用語リストを戻すステップであって、前記複数の提案された用語リストは、前記用語／フレーズが1以上の用語クラスタ内の用語に一致するときに、生成されるステップと

を含むことを特徴とするコンピュータ実行可能方法。

【請求項2】

前記エンティティは、コンピュータプログラムアプリケーションおよび／またはエンドユーザであることを特徴とする請求項1に記載の方法。

【請求項3】

クエリログから履歴クエリを収集するステップと、

前記履歴クエリの中の高 FOO を伴う履歴クエリを決定するステップと

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4】

前記用語クラスタを作成する前に、

前記用語ベクトルのサイズを縮小するステップと、

前記用語ベクトルを正規化するステップと

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5】

評価するステップは、

前記用語／フレーズと1つまたは複数の用語クラスタからの1つ（または複数）の用語／フレーズとの間のマッチを識別するステップと、

識別するステップに応答して、前記1つ（または複数）の用語／フレーズを含む1つ（または複数）の関連語提案を生成するステップと

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6】

前記1つ（または複数）の関連語提案は、前記1つ（または複数）の用語／フレーズのうちの用語／フレーズごとに、前記用語／フレーズが1組の調べられた履歴クエリ中で発生する回数を示す出現度数値をさらに含むことを特徴とする請求項5に記載の方法。

【請求項7】

10

20

30

40

50

前記用語クラスタを生成するステップは、
 前記各高 F O O 履歴クエリを前記検索エンジンに送信して前記検索結果を取得するステップと、
 前記各高 F O O 履歴クエリに対応する少なくとも 1 サブセットの検索結果からフィーチャを抽出するステップと、
 前記フィーチャから用語頻度および逆用語頻度に応じて用語ベクトルを生成するステップと
 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 8】

前記フィーチャは、前記高 F O O 履歴クエリ用語ごとにタイトル、簡単な説明、および / またはコンテキストを含むことを特徴とする請求項 7 に記載の方法。

10

【請求項 9】

前記フィーチャを抽出するステップの検索結果のサブセットは、前記検索結果のうちの上位順位付けされた結果を含むことを特徴とする請求項 7 に記載の方法。

【請求項 10】

前記用語クラスタは、第 1 の組の用語クラスタであり、前記方法は、
 前記用語 / フレーズと前記複数の用語 / フレーズとの間にマッチが存在しないことを判定するステップと、

前記判定するステップに回答して、
 各用語ベクトルが、前記検索エンジンに以前にサブMITTされた 1 組の低 F O O 履歴クエリに関連する検索結果から生成される用語ベクトルの計算された類似性から第 2 の組の用語クラスタを作成するステップと、

20

1 つまたは複数の関連語提案を識別するために前記第 2 の組の用語クラスタのうちの用語 / フレーズを考慮して前記用語 / フレーズを評価するステップと

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】

作成するステップは、
 クエリログから調べられる履歴クエリから前記低 F O O 履歴クエリを識別するステップと、

少なくとも 1 サブセットの前記各低 F O O 履歴クエリを前記検索エンジンに送信して検索結果を取得するステップと、

30

少なくとも 1 サブセットの検索結果からフィーチャを抽出するステップと、
 前記フィーチャから用語頻度および逆用語頻度に応じて前記用語ベクトルを生成するステップと

をさらに含むことを特徴とする請求項 10 に記載の方法。

【請求項 12】

クラスタリングの後に、
 前記用語 / フレーズと高 F O O 履歴クエリに基づいた前記第 1 の組の用語クラスタからの 1 つ (または複数) の用語 / フレーズとの間にマッチが存在しないことを判定するステップと、

40

前記判定するステップに回答して、前記用語 / フレーズと低 F O O 履歴クエリに基づいた前記第 2 の組の 1 つまたは複数の用語クラスタからの 1 つ (または複数) の用語 / フレーズとの間のマッチを識別するステップと、

識別するステップに回答して、前記 1 つ (または複数) の用語 / フレーズを含む 1 つ (または複数) の関連語提案を生成するステップと

をさらに含むことを特徴とする請求項 11 に記載の方法。

【請求項 13】

プロセッサに、以下の方法を実行させるためのコンピュータ実行可能命令を格納したコンピュータ読取り可能記憶媒体であって、前記方法は、

マルチセンスクエリを介して、検索結果を調べるステップであって、

50

前記マルチセンスクエリが、以下のステップ、すなわち、
 サブミットした用語／フレーズに意味的に関連している用語／フレーズを判定するステップであって、意味関係が、意味を判定するために、用語／フレーズの内容を検索することにより、発見されるステップと、
 出現度数（F O O）のしきい値を構成するステップと、
 前記構成されたしきい値に基づき、高F O O、または、低F O Oを履歴クエリに割り振るステップと、
 検索エンジンに以前にサブミットされた1組の高出現度数（F O O）履歴クエリに関連する検索結果から用語ベクトルを生成するステップと、
 用語ベクトルの計算された類似性に応じて用語クラスタを生成するステップであって、
 前記計算された類似性である $sim(q_j, q_k)$ は、以下の式、すなわち、

10

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

に従って決定されるステップであって、
 式中で、 d は、ベクトル次数を表し、 q は、クエリを表し、 k は、次数インデックスであり、

20

i 番目のベクトルの j 番目の用語についての重み w は、以下の式、すなわち

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

に従って計算され、

式中で、 TF_{ij} は、用語頻度を表し、 N は、クエリ用語の合計数であり、 DF_j は、 i 番目のベクトルの j 番目の用語を含む抽出済みのフィーチャレコード数である前記決定されるステップと

を含む前記検索結果を調べるステップと、

エンティティから用語／フレーズを受信したことに応答して、1つまたは複数の関連語提案を識別するために、前記用語クラスタ中の用語／フレーズを考慮して、前記マルチセンスクエリを介して前記用語／フレーズを評価するステップであって、前記識別することは、F O Oと確信度値との組み合わせに基づいて行われることを特徴とする評価するステップと、

30

前記F O Oと確信度値との組み合わせにより順序付けられる少なくとも1つの提案された用語リストを戻すステップであって、前記複数の提案された用語リストは、前記用語／フレーズが1以上の用語クラスタ内の用語に一致するときに、生成されるステップと

を含むことを特徴とするコンピュータ読取り可能記録媒体。

【請求項14】

前記エンティティは、コンピュータプログラムアプリケーションおよび／またはエンドユーザであることを特徴とする請求項13に記載のコンピュータ読取り可能記録媒体。

40

【請求項15】

クエリログから履歴クエリ用語を収集するコンピュータ実行可能命令と、
 前記履歴クエリ用語のうちの高F O Oを伴う履歴クエリ用語を決定するコンピュータ実行可能命令と

をさらに含むことを特徴とする請求項13に記載のコンピュータ読取り可能記録媒体。

【請求項16】

前記用語クラスタを作成する前に、
 前記用語ベクトルのサイズを縮小するコンピュータ実行可能命令と、
 前記用語ベクトルを正規化するコンピュータ実行可能命令と
 をさらに含むことを特徴とする請求項13に記載のコンピュータ読取り可能記録媒体。

50

【請求項 17】

マルチセンスクエリを介して、評価するステップは、

前記用語／フレーズと1つまたは複数の用語クラスタからの1つ（または複数）の用語／フレーズとの間のマッチを識別するコンピュータ実行可能命令と、

識別するステップに回答して、前記1つ（または複数）の用語／フレーズを含む1つ（または複数）の関連語提案を生成するコンピュータ実行可能命令と

をさらに含むことを特徴とする請求項 13 に記載のコンピュータ読取り可能記録媒体。

【請求項 18】

前記1つ（または複数）の関連語提案は、前記1つ（または複数）の用語／フレーズのうちの用語／フレーズごとに、前記用語／フレーズが1組の調べられた履歴クエリ中で発生する回数を示す出現度数値をさらに含むことを特徴とする請求項 17 に記載のコンピュータ読取り可能記録媒体。

10

【請求項 19】

前記用語クラスタを生成するステップは、

前記各高 F O O 履歴クエリを前記検索エンジンに送信して前記検索結果を取得するコンピュータ実行可能命令と、

前記各高 F O O 履歴クエリに対応する少なくとも1サブセットの検索結果からフィーチャを抽出するコンピュータ実行可能命令と、

前記フィーチャから用語頻度および逆用語頻度に応じて用語ベクトルを生成するコンピュータ実行可能命令と

20

をさらに含むことを特徴とする請求項 13 に記載のコンピュータ読取り可能記録媒体。

【請求項 20】

前記フィーチャは、前記高 F O O 履歴クエリ用語ごとにタイトル、簡単な説明、および／またはコンテキストを含むことを特徴とする請求項 19 に記載のコンピュータ読取り可能記録媒体。

【請求項 21】

前記フィーチャを抽出するステップの検索結果のサブセットは、前記検索結果のうちの上位順位付けされた結果を含むことを特徴とする請求項 19 に記載のコンピュータ読取り可能記録媒体。

【請求項 22】

30

前記用語クラスタは、第1の組の用語クラスタであり、前記コンピュータ実行可能命令は、

前記用語／フレーズと前記複数の用語／フレーズとの間にマッチが存在しないことを判定する命令と、

前記判定するステップに回答して、

各用語ベクトルが、前記検索エンジンに以前にサブMITTされた1組の低 F O O 履歴クエリに関連する検索結果から生成される用語ベクトルの計算された類似性から第2の組の用語クラスタを作成する命令と、

前記第2の組の用語クラスタのうちの用語／フレーズを考慮して前記用語／フレーズを評価し、1つまたは複数の関連語提案を識別する命令と

40

をさらに含むことを特徴とする請求項 13 に記載のコンピュータ読取り可能記録媒体。

【請求項 23】

第2の組の用語クラスタを作成するステップは、

クエリログから調べられる履歴クエリから前記低 F O O 履歴クエリを識別するコンピュータ実行可能命令と、

少なくとも1サブセットの前記各低 F O O 履歴クエリを前記検索エンジンに送信して検索結果を取得するコンピュータ実行可能命令と、

少なくとも1サブセットの検索結果からフィーチャを抽出するコンピュータ実行可能命令と、

前記フィーチャから用語頻度および逆用語頻度に応じて前記用語ベクトルを生成するコ

50

コンピュータ実行可能命令と

をさらに含むことを特徴とする請求項 2 2 に記載のコンピュータ読取り可能記録媒体。

【請求項 2 4】

クラスタリングの後に、

前記用語 / フレーズと高 F O O 履歴クエリに基づいた前記第 1 の組の用語クラスタからの 1 つ (または複数) の用語 / フレーズとの間にマッチが存在しないことを判定するコンピュータ実行可能命令と、

前記判定するステップにตอบสนองして、前記用語 / フレーズと低 F O O 履歴クエリに基づいた前記第 2 の組の 1 つまたは複数の用語クラスタからの 1 つ (または複数) の用語 / フレーズとの間のマッチを識別するコンピュータ実行可能命令と、

識別するステップにตอบสนองして、前記 1 つ (または複数) の用語 / フレーズを含む 1 つ (または複数) の関連語提案を生成するコンピュータ実行可能命令と

をさらに含むことを特徴とする請求項 2 3 に記載のコンピュータ読取り可能記録媒体。

【請求項 2 5】

プロセッサと、

前記プロセッサに結合されたメモリであって、前記プロセッサに、以下の方法を実行させるためのコンピュータ実行可能命令を格納した前記メモリとを備え、

前記方法は、

マルチセンスクエリを介して、検索結果を調べるステップであって、

前記マルチセンスクエリが、以下のステップ、すなわち、

サブミットした用語 / フレーズに意味的に関連している用語 / フレーズを判定するステップであって、意味関係が、意味を判定するために、用語 / フレーズの内容を検索することにより、発見されるステップと、

出現度数 (F O O) のしきい値を構成するステップと、

前記構成されたしきい値に基づき、高 F O O、または、低 F O O を履歴クエリに割り振るステップと、

検索エンジンに以前にサブミットされた 1 組の高出現度数 (F O O) 履歴クエリに関連する検索結果から用語ベクトルを生成するステップと、

用語ベクトルの計算された類似性に応じて用語クラスタを生成するステップであって、前記計算された類似性である $sim(q_j, q_k)$ は、以下の式、すなわち、

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

に従って決定されるステップであって、

式中で、 d は、ベクトル次数を表し、 q は、クエリを表し、 k は、次数インデックスであり、

i 番目のベクトルの j 番目の用語についての重み w は、以下の式、すなわち

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

に従って計算され、

式中で、 TF_{ij} は、用語頻度を表し、 N は、クエリ用語の合計数であり、 DF_j は、 i 番目のベクトルの j 番目の用語を含む抽出済みのフィーチャレコード数である前記決定されるステップと

を含む前記検索結果を調べるステップと、

エンティティから用語 / フレーズを受信したことにตอบสนองして、1 つまたは複数の関連語提案を識別するために、前記用語クラスタ中の用語 / フレーズを考慮して、前記マルチセンスクエリを介して前記用語 / フレーズを評価するステップであって、前記識別すること

10

20

30

40

50

は、F O Oと確信度値との組み合わせに基づいて行われることを特徴とする評価するステップと、

前記F O Oと確信度値との組み合わせにより順序付けられる少なくとも1つの提案された用語リストを戻すステップであって、前記複数の提案された用語リストは、前記用語/フレーズが1以上の用語クラスタ内の用語に一致するときに、生成されるステップとを含む、

ことを特徴とするコンピューティングデバイス。

【請求項26】

前記エンティティは、コンピュータプログラムアプリケーションおよび/またはエンドユーザであることを特徴とする請求項25に記載のコンピューティングデバイス。

10

【請求項27】

クエリログから履歴クエリ用語を収集するコンピュータ実行可能命令と、
前記履歴クエリ用語のうちの高F O Oを伴う履歴クエリ用語を決定するコンピュータ実行可能命令と

をさらに含むことを特徴とする請求項25に記載のコンピューティングデバイス。

【請求項28】

前記用語クラスタを作成する前に、

前記用語ベクトルのサイズを縮小するコンピュータ実行可能命令と、

前記用語ベクトルを正規化するコンピュータ実行可能命令と

をさらに含むことを特徴とする請求項25に記載のコンピューティングデバイス。

20

【請求項29】

評価するステップは、

前記用語/フレーズと1つまたは複数の用語クラスタからの1つ(または複数)の用語/フレーズとの間のマッチを識別するコンピュータ実行可能命令と、

識別するステップにตอบสนองして、前記1つ(または複数)の用語/フレーズを含む1つ(または複数)の関連語提案を生成するコンピュータ実行可能命令と

をさらに含むことを特徴とする請求項25に記載のコンピューティングデバイス。

【請求項30】

前記1つ(または複数)の関連語提案は、前記1つ(または複数)の用語/フレーズのうちの用語/フレーズごとに、前記用語/フレーズが1組の調べられた履歴クエリ中で発生する回数を示す出現度数値をさらに含むことを特徴とする請求項29に記載のコンピューティングデバイス。

30

【請求項31】

前記用語クラスタを生成するステップは、

前記各高F O O履歴クエリを前記検索エンジンに送信して前記検索結果を取得するコンピュータ実行可能命令と、

前記各高F O O履歴クエリに対応する少なくとも1サブセットの検索結果からフィーチャを抽出するコンピュータ実行可能命令と、

前記フィーチャから用語頻度および逆用語頻度に応じて用語ベクトルを生成するコンピュータ実行可能命令と

をさらに含むことを特徴とする請求項25に記載のコンピューティングデバイス。

40

【請求項32】

前記フィーチャは、前記高F O O履歴クエリ用語ごとにタイトル、簡単な説明、および/またはコンテキストを含むことを特徴とする請求項31に記載のコンピューティングデバイス。

【請求項33】

前記フィーチャを抽出するステップの検索結果のサブセットは、前記検索結果のうちの上位順位付けされた結果を含むことを特徴とする請求項31に記載のコンピューティングデバイス。

【請求項34】

50

前記用語クラスタは、第 1 の組の用語クラスタであり、前記コンピュータ実行可能命令は、

前記用語 / フレーズと前記複数の用語 / フレーズとの間にマッチが存在しないことを判定する命令と、

前記判定するステップにตอบสนองして、

各用語ベクトルが、前記検索エンジンに以前にサブMITされた 1 組の低 F O O 履歴クエリに関連する検索結果から生成される用語ベクトルの計算された類似性から第 2 の組の用語クラスタを作成する命令と、

前記第 2 の組の用語クラスタのうちの用語 / フレーズを考慮して前記用語 / フレーズを評価し、1 つまたは複数の関連語提案を識別する命令と

をさらに含むことを特徴とする請求項 2 5 に記載のコンピューティングデバイス。

【請求項 3 5】

第 2 の組の用語クラスタを作成するステップは、

クエリログから調べられる履歴クエリから前記低 F O O 履歴クエリを識別するコンピュータ実行可能命令と、

少なくとも 1 サブセットの前記各低 F O O 履歴クエリを前記検索エンジンに送信して検索結果を取得するコンピュータ実行可能命令と、

少なくとも 1 サブセットの検索結果からフィーチャを抽出するコンピュータ実行可能命令と、

前記フィーチャから用語頻度および逆用語頻度に応じて前記用語ベクトルを生成するコンピュータ実行可能命令と

をさらに含むことを特徴とする請求項 3 4 に記載のコンピューティングデバイス。

【請求項 3 6】

クラスタリングの後に、

前記用語 / フレーズと高 F O O 履歴クエリに基づいた前記第 1 の組の用語クラスタからの 1 つ (または複数) の用語 / フレーズとの間にマッチが存在しないことを判定するコンピュータ実行可能命令と、

前記判定するステップにตอบสนองして、前記用語 / フレーズと低 F O O 履歴クエリに基づいた前記第 2 の組の 1 つまたは複数の用語クラスタからの 1 つ (または複数) の用語 / フレーズとの間のマッチを識別するコンピュータ実行可能命令と、

識別するステップにตอบสนองして、前記 1 つ (または複数) の用語 / フレーズを含む 1 つ (または複数) の関連語提案を生成するコンピュータ実行可能命令と

をさらに含むことを特徴とする請求項 3 5 に記載のコンピューティングデバイス。

【請求項 3 7】

プロセッサと、

マルチセンスクエリを介して、検索結果を調べる手段であって、

前記マルチセンスクエリが、以下の手段、すなわち、

サブMITした用語 / フレーズに意味的に関連している用語 / フレーズを判定するステップであって、意味関係が、意味を判定するために、用語 / フレーズの内容を検索することにより、発見される手段と、

出現度数 (F O O) のしきい値を構成する手段と、

前記構成されたしきい値に基づき、高 F O O、または、低 F O O を履歴クエリに割り振る手段と、

検索エンジンに以前にサブMITされた 1 組の高出現度数 (F O O) 履歴クエリに関連する検索結果から用語ベクトルを生成する手段と、

用語ベクトルの計算された類似性に応じて用語クラスタを生成する手段であって、前記計算された類似性である $s i m (q_j, q_k)$ は、以下の式、すなわち、

10

20

30

40

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

に従って決定される手段であって、

式中で、 d は、ベクトル次数を表し、 q は、クエリを表し、 k は、次数インデックスであり、

i 番目のベクトルの j 番目の用語についての重み w は、以下の式、すなわち

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

に従って計算され、

式中で、 TF_{ij} は、用語頻度を表し、 N は、クエリ用語の合計数であり、 DF_j は、 i 番目のベクトルの j 番目の用語を含む抽出済みのフィーチャレコード数である前記決定される手段と

を含む前記検索結果を調べる手段と、

エンティティから用語／フレーズを受信したことに応答して、1つまたは複数の関連語提案を識別するために、前記用語クラスタ中の用語／フレーズを考慮して、前記マルチセンスクエリを介して前記用語／フレーズを評価する手段であって、前記識別することは、 FOO と確信度値との組み合わせに基づいて行われることを特徴とする評価する手段と、

前記 FOO と確信度値との組み合わせにより順序付けられる少なくとも1つの提案された用語リストを戻す手段であって、前記複数の提案された用語リストは、前記用語／フレーズが1以上の用語クラスタ内の用語に一致するときに、生成される手段と

を含む、

ことを特徴とするコンピューティングデバイス。

【請求項38】

前記エンティティは、コンピュータプログラムアプリケーションおよび／またはエンドユーザであることを特徴とする請求項37に記載のコンピューティングデバイス。

【請求項39】

クエリログから履歴クエリ用語を収集する収集手段と、

前記履歴クエリ用語のうちの高 FOO を伴う履歴クエリ用語を決定する決定手段と

をさらに備えることを特徴とする請求項37に記載のコンピューティングデバイス。

【請求項40】

前記評価手段は、

前記用語／フレーズと1つまたは複数の用語クラスタからの1つ（または複数）の用語／フレーズとの間のマッチを識別する識別手段と、

識別するステップに応答して、前記1つ（または複数）の用語／フレーズを含む1つ（または複数）の関連語提案を生成する生成手段と

をさらに備えることを特徴とする請求項37に記載のコンピューティングデバイス。

【請求項41】

前記用語クラスタを生成する前記生成手段は、

前記各高 FOO 履歴クエリを前記検索エンジンに送信して前記検索結果を取得する送信手段と、

前記各高 FOO 履歴クエリに対応する少なくとも1サブセットの検索結果からフィーチャを抽出する抽出手段と、

前記フィーチャから用語ベクトルを生成する生成手段と

をさらに備えることを特徴とする請求項37に記載のコンピューティングデバイス。

【請求項42】

前記用語クラスタは、第1の組の用語クラスタであり、前記コンピューティングデバイ

10

20

30

40

50

スは、

前記用語／フレーズと前記複数の用語／フレーズとの間にマッチが存在しないことを判定する判定手段と、

前記判定するステップに応答して、

各用語ベクトルが、前記検索エンジンに以前にサブミットされた1組の低F O O履歴クエリに関連する検索結果から生成される用語ベクトルの計算された類似性から第2の組の用語クラスタを作成する作成手段と、

前記第2の組の用語クラスタのうちの用語／フレーズを考慮して前記用語／フレーズを評価し、1つまたは複数の関連語提案を識別する評価手段と

をさらに備えることを特徴とする請求項37に記載のコンピューティングデバイス。 10

【請求項43】

前記用語／フレーズと高F O O履歴クエリに基づいた前記第1の組の用語クラスタからの1つ（または複数）の用語／フレーズとの間にマッチが存在しないことを計算する計算手段と、

前記計算するステップに応答して、前記用語／フレーズと低F O O履歴クエリに基づいた前記第2の組の1つまたは複数の用語クラスタからの1つ（または複数）の用語／フレーズとの間のマッチを識別する識別手段と、

識別するステップに応答して、前記1つ（または複数）の用語／フレーズを含む1つ（または複数）の関連語提案を生成する生成手段と

をさらに備えることを特徴とする請求項42に記載のコンピューティングデバイス。 20

【発明の詳細な説明】

【技術分野】

【0001】

本発明のシステムおよび方法は、データマイニングに関する。

【背景技術】

【0002】

キーワードまたはフレーズは、WWW（ワールドワイドウェブ）上の関連したウェブページ／サイトを求めて検索する際におけるウェブサーファ（Web surfer）が検索エンジンにサブミット（submit）する単語または1組の用語である。検索エンジンは、そのページ／サイト上に現れるこれらのキーワードおよびキーワードフレーズに基づいてウェブサイトの関連性を決定する。ウェブサイトトラフィックのかなりの割合は、検索エンジンの使用からもたらされるので、ウェブサイトプロモータは、適切なキーワード／フレーズの選択が、サイトトラフィックを増加させて所望のサイトを見せるために不可欠であることを知っている。検索エンジン結果最適化（search engine result optimization）のためにウェブサイトに関連するキーワードを識別する技法は、例えばウェブサイトのコンテンツ、および関連したキーワードを識別する目的の人間による評価を含んでいる。この評価には、キーワード人気ツール（keyword popularity tool）の使用を含むこともある。かかるツールは、どれだけ多くの人々が、検索エンジンに対するキーワードを含めて、特定のキーワードまたはフレーズをサブミットしたかを決定する。そのウェブサイトに関連しており、検索クエリを生成するに際して非常に頻繁に使用されると決定されたキーワードが、一般にそのウェブサイトに関して検索エンジン結果最適化のために選択される。 30

【0003】

そのウェブサイトの検索エンジン結果最適化のために1組のキーワードを識別した後に、プロモータは、その検索エンジン結果中で（他のウェブサイトの検索エンジン結果の表示位置に比べて）さらに高い位置にウェブサイトを進めたいと望むこともある。この目的を達成するために、このプロモータは、この1つ（または複数）のキーワードに値を付けて、ウェブサーファがこの1つ（または複数）のキーワードに関連するこのプロモータのリストをクリックするたびに、このプロモータがいくら支払うかを指し示す。換言すれば、キーワードの入札が、ペイパークリック入札（pay-per-click bid） 40 50

になっている。この同じキーワードについての他の入札に比べてこのキーワード入札の金額が多額になればなるほど、この検索エンジンは、このキーワードに基づいて検索結果中でより高い位置に（重要性に関してより目立つように）この関連するウェブサイトを表示することになる。

【発明の開示】

【発明が解決しようとする課題】

【0004】

以上に鑑みて、ウェブサイトコンテンツに関連したキーワードをより良好に識別するシステムおよび方法があれば、ウェブサイトプロモータによって歓迎されるはずである。これにより、これらのプロモータは、ユーザの好ましい用語を入札することができるようになる。これらのシステムおよび方法が、人間がウェブサイトコンテンツを評価して、検索エンジン最適化およびキーワード入札のための関連したキーワードを識別する必要性と無関係になるとすれば、それは理想的である。

【課題を解決するための手段】

【0005】

関連語提案 (related term suggestion) についてのシステムおよび方法について説明している。一態様においては、用語クラスタ (term cluster) が、用語ベクトル (term vector) の計算された類似性に依拠して生成される。検索結果から生成されている各用語ベクトルは、検索エンジンに以前にサブミットされた1組の高FOO (frequency of occurrence 出現度数) 履歴クエリに関連したものである。あるエンティティから用語/フレーズを受信するのに応答して、この用語/フレーズを、用語クラスタ中の用語/フレーズに鑑みて評価して、1つまたは複数の関連語提案を識別する。

【0006】

図面中において、コンポーネント参照番号のいちばん左側の桁は、そのコンポーネントが最初に現れる特定の図面を識別している。

【発明を実施するための最良の形態】

【0007】

概要

関連した用語/フレーズを提案する最も簡単な方法は、サブstring マッチングアプローチ (substring matching approach) を使用することであると思われ、このアプローチは、1つの用語/フレーズが、別の用語/フレーズの単語の一部または全部を含むときに2つの用語/フレーズが関連していると判断するものである。しかし、この技法は、実質的には制限されている。関連語は、共通の単語を含む必要がないので、この方法は、多くの意味的に関連した用語を無視してしまうこともある。例えば、履物会社が「shoe (靴)」についての関連語を知りたいと思うことについて考察してみる。従来のマッチングアプローチが使用される場合には、「women's shoes (婦人靴)」、「discount shoes (ディスカウント靴)」などだけしか提案されないことになる。しかし、「sneakers (スニーカー)」、「hiking boots (ハイキングブーツ)」、「Nike (ナイキ)」などの他の多くの関連語も存在する。

【0008】

マルチセンスクエリについての関連語を提案するための以下のシステムおよび方法は、従来のサブstring マッチング技法のこれらの制限に対処している。この目的を達成するためにこれらのシステムおよび方法では、エンドユーザ（例えば、ウェブサイトプロモータ、広告主など）がサブミットした用語/フレーズに意味的に関連している用語/フレーズについての情報を得るために検索エンジン結果が調べられる。この意味的な関係については、用語/フレーズの意味を考慮することができる用語/フレーズを取り巻くコンテキスト（例えば、テキストなど）をこれらの検索エンジン結果から調べることによって構築することができる。より詳細には、1組のクエリ用語を履歴クエリログから集めて、そ

10

20

30

40

50

これらの出現度数 (F O O) がカウントされる。これらのクエリ用語は、1つずつその検索エンジンにサブMITTされる。一実装形態においては、このサブMITTされた履歴クエリログ用語は、他の履歴クエリログ用語の出現度数に比べて相対的に高い出現度数を有している。

【 0 0 0 9 】

これらのサブMITTされた各クエリを受信するのに応答して、この検索エンジンは、このサブMITTされたクエリの周囲の U R L、結果タイトル、各結果の簡単な説明および/またはコンテキストを含めて、検索結果の順位付けされたリストを返す。検索エンジン結果が受信された後に、これらのシステムおよび方法は、これらの返された検索結果 (例えば、1つまたは複数の上位に順位付けされた結果) のうちの選択された結果から1組のフィーチャ (キーワードおよび知られている T F I D F 技法を使用して計算される対応する重み) を抽出する。これらのサブMITTされた検索クエリからの対応する検索エンジン結果のフィーチャを抽出した後に、これらの抽出されたフィーチャが正規化される。これらの正規化されたフィーチャを使用して、各サブMITTされたクエリが表され、またこれらの正規化されたフィーチャをテキストクラスタリングアルゴリズム (t e x t c l u s t e r i n g a l g o r i t h m) 中で使用して、サブMITTされたクエリ用語がクラスタにグループ化される。

10

【 0 0 1 0 】

そのエンドユーザからこの用語/フレーズを受信するのに応答して、この用語/フレーズは、これらの用語クラスタ中の各用語/フレーズと比較される。これらの用語クラスタは互いにコンテキスト的に関連した用語を含んでいるので、この用語/フレーズがこれらのクラスタ内の用語と比較されるときに、この用語/フレーズは、任意の複数の関連したコンテキストまたは「意味」に鑑みて評価される。一実装形態においては、用語/フレーズがクラスタからの用語とマッチングする場合、このクラスタは、提案される用語リストの形でそのエンドユーザに返される。この提案される用語リストは、この用語/フレーズと、各用語/フレーズ間の類似性評価値 (確信度値) と、各用語/フレーズの出現度数 (F O O) とに意味的および/またはコンテキスト的に関連づけて決定される用語/フレーズを含んでいる。この返されたリストは、 F O O および確信度値の組合せによって順序付けられる。この用語/フレーズが、複数の用語クラスタ中の用語とマッチングする場合には、複数の提案される用語リストが生成される。これらのリストは、これらのクラスタサイズによって順序付けられる。また各リスト内の用語は、 F O O および確信度値の組合せによって順序付けられる。マッチングクラスタが識別されない場合には、このクエリ用語はさらに、低 F O O のクエリ用語から生成される拡張クラスタに対してマッチングが行われる。

20

30

【 0 0 1 1 】

一実装形態においては、低 F O O を伴うクエリ用語は、高出現度数の履歴クエリログ用語から生成される用語クラスタについて分類機構 (c l a s s i f i e r) (例えば、K個の最近接分類機構) をトレーニングすることによってクラスタ化される。低出現度数を有すると決定された履歴クエリ用語が、1つずつこの検索エンジンに対してサブMITTされる。次いで、フィーチャが、この返された検索結果のうちの選択された結果 (例えば、第1の上位に順位付けされたウェブページなど) から抽出される。この抽出されたフィーチャは正規化され、これを使用して低 F O O を伴うクエリ用語が表される。次いで、これらのクエリ用語を既存のクラスタに分類して、このトレーニングされた分類機構に基づいて拡張クラスタが生成される。次いでこのエンドユーザがサブMITTした用語/フレーズをこれらの拡張クラスタに鑑みて評価して、提案される用語リストを識別しこのエンドユーザに返す。

40

【 0 0 1 2 】

次に、マルチセンスクエリについての関連した用語/キーワードを提案するためのこれらおよび他の態様のシステムおよび方法について、より詳細に説明する。

【 0 0 1 3 】

50

例示のシステム

図面を参照すると、マルチセンスクエリについての関連語を提案するためのシステムおよび方法が、適切なコンピューティング環境において実装されるものとして説明され示されている。図面では、同様な参照番号は、同様なエレメントを示している。必要という訳ではないが、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される。プログラムモジュールは一般に、特定のタスクを実施し、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、データ構造などを含んでいる。システムおよび方法は、前述の文脈で説明されるが、以降に説明する動作およびオペレーションは、ハードウェアで実装することもできる。

10

【0014】

図1は、マルチセンスクエリについての関連語を提案するための例示のシステム100を示している。この実装形態においては、システム100は、ネットワーク104を横切ってクライアントコンピューティングデバイス106に結合されたEVS (editorial verification server 編集検証サーバ) 102を含んでいる。例えばクライアントコンピューティングデバイス106から、またはEVS 102上で実行される別のアプリケーション（図示せず）から用語/フレーズ108を受信するのに対応して、EVS 102は、提案される用語リスト110を生成し、このクライアントコンピューティングデバイス106に伝えて、これによりエンドユーザは、この用語/フレーズに実際に入札するのに先立ってこの用語/フレーズ108に意味的/コンテキスト的に関連する1組の用語を評価できるようになる。ネットワーク104は、オフィス、企業規模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において一般的になっている通信環境など、LAN (ローカルエリアネットワーク) 通信環境、および一般的なWAN (ワイドエリアネットワーク) 通信環境の任意の組合せを含むことができる。システム100が、クライアントコンピューティングデバイス106を含むとき、このクライアントコンピューティングデバイスは、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ティングデバイス（例えば、セルラ電話、携帯情報端末 (personal digital assistant)、またはハンドヘルドコンピュータ) など、どのようなタイプのコンピューティングデバイスでもよい。

20

【0015】

提案される用語リスト110は、例えば、この用語/フレーズ108と、各用語/フレーズ108間の類似性評価値（確信度値）と、各用語/フレーズ出現度数（FOO）（履歴クエリログにおける頻度）とに関連していると決定された用語/フレーズを含んでいる。関連した用語/フレーズを識別し、類似性評価値を生成し、FOO値を生成するための技法については、キーワードマイニング (keyword mining)、フィーチャ抽出、および用語クラスタリングとタイトルが付けられたセクションに関連して以下で非常に詳細に説明している。

30

【0016】

表1は、「mail (メール)」の用語/フレーズ108に関連していると決定された用語の例示の提案される用語リスト110を示している。用語/フレーズ108に関連した用語は、この実施例においては、提案される用語とタイトルが付けられたカラム1中に示される。

40

【0017】

【表 1】

表 1

入札用語「MAIL」についての例示の提案される用語リスト

提案される用語	類似性	頻度	<コンテキスト>
hotmail	0.246142	93161	オンライン電子メール関連
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	0.0968248	4481	
free email	0.130611	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	0.0800538	108	
webmail.yahoo.com	0.08789	71	
free email account	0.0836481	65	
提案される用語	類似性	頻度	
mail	1	2191	従来メール関連
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	

【0018】

表 1 を参照すると、この提案される用語リスト中の用語は、用語類似性値（「類似性」とタイトルの付けられたカラム 2 を参照）と出現度数スコア（「頻度」とタイトルが付けられたカラム 3 を参照）にマッピングされることに留意されたい。「用語クラスタリング」とタイトルが付けられたセクション中において以下で説明しているように計算された各用語の類似性値は、対応する提案される用語（カラム 1）とこの実施例においては「mail」となる用語 / フレーズ 108 との間の類似性評価値を提供する。各頻度値またはス

10

20

30

40

50

コアは、この提案される用語が履歴クエリログ中に現れる回数を指し示す。この提案される用語リストは、用語の類似性、および/またはビジネスゴールに応じた出現度数スコアに応じてソートされる。

【0019】

所与の任意の用語/フレーズ108(例えば、mailなど)は、その内部で入札用語を使用することができる複数のコンテキストを有することもある。これを明らかにするために、STSモジュール112は、どの提案される用語の提案される用語リスト110が用語/フレーズ108の複数のコンテキストのうちのどれに対応するかについての指示を提供する。例えば、表1を参照すると、「mail」の用語/フレーズ108は、2つのコンテキスト、すなわち(1)従来のオフラインメール、および(2)オンライン電子メールを有する。関連語の各リストが、これらの2つの各入札用語コンテキストについて示されることに留意されたい。

10

【0020】

さらに、任意の用語/フレーズ108についての提案される用語は、入札用語の同義語以上のものとすることができる。例えば、表1を参照すると、この提案される用語「usps」は、メールを扱う組織についての頭字語(acronym)であるが、入札用語「mail(メール)」についての同義語ではない。しかし、「usps」は、「mail」入札用語に非常に関連した用語でもあり、したがってこの提案される用語リスト110中に示される。一実装形態においては、STSモジュール112は、関連付けルール、すなわちitr(T) itr(R)に応じて関連語R(例えば「usps」とターゲット用語T(例えば「mail」)の関係性を決定し、式中で「itr」は、「関心がある(interested in)」を表している。ユーザ(広告主、および/またはウェブサイトプロモータなど)が、Rに関心がある場合、このユーザは、Tにも関心があることになる。

20

【0021】

EVS102は、提案される用語リスト110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んでいる。このコンピュータプログラムモジュールは、例えば、STS(search term suggestion検索用語提案)モジュール112および分類モジュール114を含んでいる。STSモジュール112は、クエリログ118から1組の履歴クエリ116を検索する。この履歴クエリは、検索エンジンに対して以前にサブMITTされた検索クエリ用語を含んでいる。STSモジュール112は、出現度数に応じて履歴クエリ116を評価して、高出現度数(FOO)検索用語120と相対的により低い出現度数検索用語122を識別する。この実装形態においては、構成可能なしきい値を使用して、履歴クエリが比較的より高い出現度数を有するか、それとも低い出現度数を有するかを判定する。例えば、少なくともしきい値回数だけ現れる履歴クエリ116中の検索クエリ用語は、高出現度数を有していると言われる。同様に、しきい値回数よりも少なくしか現れない履歴クエリ116中の検索クエリ用語は、低出現度数を有していると言われる。図では、かかるしきい値は、「他のデータ」124の各部分として示される。

30

【0022】

キーワードマイニングおよびフィーチャ抽出

STSモジュール112は、各クエリ(検索クエリ128)を1つずつ検索エンジン126にサブMITTすることにより、意味的/コンテキスト的な意味の高出現度数クエリ用語120を調べる。検索クエリ128を受信するのに応答して、検索エンジン126は、検索結果130中の(番号を構成可能な)順位付けされたリストをSTSモジュール112に返す。この順位付けされたリストは、サブMITTされた検索クエリ128に関連したURL、結果タイトル、簡単な説明および/またはクエリ用語のコンテキストを含んでいる。この順位付けされたリストは、検索結果132に記憶される。かかる検索結果の取出しは検索クエリ128ごとに行われる。

40

【0023】

50

S T Sモジュール112は、ウェブページのHTML（ハイパーテキストマークアップ言語）を解析して各取り出された検索結果132からクエリ用語120ごとにそのURL、結果タイトル、簡単な説明および/またはそのクエリ用語のコンテキストを抽出する。このURL、結果タイトル、簡単な説明および/またはそのクエリ用語のコンテキスト、ならびにこの取り出された検索結果132を取得するために使用される検索クエリ128は、抽出されたフィーチャ134の各レコードにS T Sモジュール112によって記憶される。

【0024】

高出現度数クエリ用語120についての検索結果130を解析した後に、S T Sモジュール112は、抽出されたフィーチャ134上でテキスト前処理オペレーションを実施して、この抽出されたフィーチャから言語トークンを生成（トークン化）して個別のキーワードに入れる。これらのトークンのサイズを縮小するために、S T Sモジュール112は、どのようなストップ語句（stop-word）（例えば、「the」、「a」、「is」など）も除去し、共通の接尾辞を除去して、例えば知られているポーターステミングアルゴリズム（Porter stemming algorithm）を使用してこれらのキーワードを正規化する。S T Sモジュール112は、この結果得られる抽出されたフィーチャ134を1つまたは複数の用語ベクトル136に構成する。

10

【0025】

各用語ベクトル136は、TFIDF（term frequency and inverted document frequency用語頻度および逆ドキュメント頻度）スコアに基づいたサイズを有する。i番目のベクトルのj番目のキーワードについての重みは、次式のように計算される。

20

【0026】

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

式中で、 TF_{ij} は用語頻度（i番目のレコードにおけるキーワードjの出現回数）を表し、Nは、クエリ用語の合計数であり、 DF_j は、キーワードjを含むレコード数である。

【0027】

用語クラスタリング

S T Sモジュール112は、類似の用語をグループ化して用語ベクトル136から用語クラスタ138を生成する。この目的を達成するために、この実装形態においては、各用語のベクトル表現を仮定して、コサインファンクションを使用して1対の用語の間の類似性を評価する（これらのベクトルが正規化されていることを想起されたい）。

30

【0028】

【数1】

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

【0029】

したがって、2つの用語間の距離（類似性評価値）は、次式のように定義される。

40

【0030】

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

かかる類似性評価値は、「他のデータ」124の各部分として示されている。かかる例示の類似性値は、表1の例示の提案される用語リスト110に示されている。

【0031】

S T Sモジュール112は、この計算された1つ（または複数）の類似性評価値を使用して、キーワードベクトル134によって表される用語をクラスタ/グループ化して1つ（または複数）の用語クラスタ138に入れる。より詳細には、この実装形態において、S T Sモジュール112は、知られているDBSCAN（density-based clustering algorithm密度ベースのクラスタリングアルゴリズム）

50

を使用して、1つ(または複数)の用語クラスタ138を生成する。DBSCANは、2つのパラメータ、すなわちEpsおよびMinPtsを使用する。Epsは、クラスタ138中のポイント間の最大距離を表す。ここでは、その尾部がその起点に移動されるときに各ベクトルはそのベクトルの頭部のポイントによって表現することができるので、ポイントは、ベクトルと等価である。MinPtsは、クラスタ138中の最小ポイント数を表す。クラスタ138を生成するためには、DBSCANは、任意ポイントpを用いて開始され、EpsおよびMinPtsに関してpからの密度の影響が及ぶすべてのポイントを取り出す。pがコアポイントである場合、このプロシージャは、EpsおよびMinPtsに関するクラスタ138をもたらす。pが境界ポイントである場合には、pからの密度の影響が及ぶポイントは存在せず、DBSCANは、次のポイントを訪れることになる。

10

【0032】

用語マッチング

エンドユーザ(例えば、広告主、ウェブサイトプロモータなど)から用語/フレーズ108を受信するのに応答して、STSモジュール112は、用語/フレーズ108を用語クラスタ138中の各用語/フレーズと比較する。用語クラスタ138は、互いにコンテキスト的に関連した用語を含んでいるので、用語/フレーズ108は、複数の関連した履歴コンテキスト、または「意味」に鑑みて評価される。一実装形態においては、用語/フレーズ108がクラスタ138からの用語/フレーズとマッチングするとSTSモジュール112が判定した場合、検索用語提案モジュール112は、クラスタ138から提案される用語リスト110を生成する。この実装形態においては、マッチは、厳密なマッチでもよく、また単数形/複数形、スペルの誤り、句読点マークなどわずかな数の変形を伴うマッチでもよい。この返されたリストは、F00と確信度値の組合せによって順序づけされる。

20

【0033】

用語/フレーズ108が複数の用語クラスタ138中の用語とマッチしているとSTSモジュール112が判定する場合、検索用語提案モジュール112は、用語クラスタ138の複数のクラスタ中の用語から複数の提案される用語リスト110を生成する。このリストは、そのクラスタサイズによって順序付けられる。また、各リスト内のこれらの用語は、F00と確信度値の組合せによって順序付けられる。

30

【0034】

低F00用語の分類

高出現度数(F00)クエリ用語120から生成された用語クラスタ138が、エンドユーザの入力用語/フレーズ108に対する同じ用語を含んでいないときには、分類モジュール114は、提案される用語リスト110を生成する。この目的を達成するために、分類モジュール114は、高出現度数(F00)クエリログ用語120から生成された用語クラスタ138からトレーニングされた分類機構140を生成する。用語クラスタ138中の用語は、すでに分類オペレーションにとって適切なベクトル空間モデル中の対応するキーワードベクトルを有している。さらにストップ語句除去、およびワードステミング(word stemming)(接尾辞除去)が、(クラスタ138が基づいている)用語ベクトル136のサイズを縮小している。一実装形態においては、追加のサイズ縮小技法、例えばフィーチャ選択または再パラメータ化を使用することができる。

40

【0035】

この実装形態においては、クラスが知られていないクエリ用語120を分類するために、分類モジュール114は、k-最近接分類機構アルゴリズム(k-Nearest Neighbor classifier algorithm)を使用して、それらの対応するフィーチャベクトルを用いてすべてのクラスが知られているクエリ用語120中のk個の最も類似した近接用語を見出し、この近接用語のクラスラベルの重み付けされた大多数を使用して、この新しいクエリ用語のクラスを予測する。ここで、すでに用語クラスタ138中に存在する各クエリ用語にはそれらの対応するクラスラベルと同じラベルが

50

割り当てられるが、各クラスタ138は、簡単な連続番号でラベル付けされる。これらの近接用語は、Xに対する各近接用語の類似性を使用して重み付けされ、ここで、類似性は、ユークリッド距離 (Euclidean distance)、または2つのベクトル間のコサイン値によって評価される。このコサイン類似性は、以下のように表される。

【0036】

【数2】

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

10

【0037】

式中で、Xは、テスト用語、すなわちベクトルとして表現される分類すべきクエリ用語である。D_jは、j番目のトレーニング用語である。t_iは、XとD_jによって共有される単語である。x_iは、X中のキーワードt_iの重みである。d_{ij}は、D_j中のキーワードt_iの重みである。

【0038】

【数3】

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

20

【0039】

は、Xのノルム (norm) であり、 $\|D_j\|_2$ はD_jのノルムである。したがって、このテスト用語Xのクラスラベルは、すべての近接用語のうちの重み付けされた大多数の近接用語のクラスラベルである。

【0040】

【数4】

$$\text{label}(X) = \arg \max_{l_i} \left(\sum_{\text{すべての } D_j, \text{ ただし } \text{label}(D_j)=l_i} \text{sim}(X, D_j) \right)$$

【0041】

30

別の実装形態においては、最近接分類技法以外の (例えば、回帰モデル (regression model)、ベイジアン分類機構 (Bayesian classifier)、判断木、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めて) 異なる統計的分類および機械学習技法を使用して、トレーニングされた分類機構140が生成される。

【0042】

分類モジュール114は、(各検索クエリ128を介して) 1つずつ、低出現度数 (FOO) クエリ用語122を検索エンジン126にサブミットする。特定の検索クエリ128に関連する検索結果130を受信し、すでに説明した技法を使用するのに応答して、分類モジュール114は、この検索結果130によって識別される1つまたは複数の取り出された検索結果132からフィーチャを抽出する (フィーチャ134を抽出している)。この実装形態においては、フィーチャは、第1の上位順位付けされた検索結果132から抽出される。取り出され解析された検索結果132ごとに、分類モジュール114は、抽出済みのフィーチャ134の各レコード中の以下の情報、すなわちこのクエリ用語のURL、結果タイトル、簡単な説明および/またはコンテキスト、ならびにこの取り出された検索結果132を取得するために使用される検索クエリ128を記憶する。次に、分類モジュール114は、トークン化を行い、サイズを縮小し、低FOOクエリ用語122から導き出された抽出済みのフィーチャを正規化して用語ベクトル136を生成する。次いで、分類モジュール114は、これらのクエリ用語をクラスタ化して各組のクラスタ138に入れる。このクラスタリングオペレーションは、(高FOOクエリ用語120から生成

40

50

された) トレーニングされた分類機構 140 を使用して実施される。

【0043】

分類モジュール 114 は、(低 F O O クエリ用語 122 に基づいて生成される) これらの拡張された用語クラスタに鑑みてエンドユーザがサブミットした用語 / フレーズ 108 を評価して、1つまたは複数の提案される用語リスト 110 を識別しこのエンドユーザに対して返す。かかる例示のプロシージャは、用語マッチングの項にて説明しており、以下のセクションにおいても説明する。

【0044】

例示のプロシージャ

図 2 は、マルチセンスクエリについての関連語を提案するための例示のプロシージャ 200 を示している。説明の目的で、このプロシージャのオペレーションについて、図 1 のコンポーネントに関連して考察する。(すべての参照番号は、このコンポーネントが最初に導入されている図面の番号から開始される)。ブロック 202 において、検索用語提案 (STS) モジュール 112 (図 1) は、クエリログ 120 から履歴クエリ用語 116 を収集する。STS モジュール 112 は、出現度数に応じて履歴クエリ 116 を構成する。ブロック 204 において、STS モジュール 112 は、高出現度数クエリ用語 120 を検索エンジン 126 に送信し、対応する検索結果 130 を受信する。ブロック 206 において、STS モジュール 112 は、各検索結果 130 から断片記述を抽出し、これらの断片記述 (抽出済みのフィーチャ 134) を一緒にマージして用語ベクトル 136 を形成する。各用語ベクトル 136 は、高出現度数クエリ用語 120 ごとに生成される。

【0045】

ブロック 208 において、STS モジュール 112 は、クラスタリングアルゴリズムを適用して、用語ベクトル 136 に基づいてほぼ同様な用語を用語クラスタ 138 にグループ分けする。ブロック 210 において、エンドユーザから用語 / フレーズ 108 を受信するのに応答して、STS モジュール 112 は、用語 / フレーズ 108 とほぼ同様であると判定された用語クラスタ 138 からの任意のキーワード / キーフレーズから、提案される用語リスト 110 を生成する。ブロック 212 において、STS モジュール 112 は、キーワードクラスタ 138 からのどのようなキーワード / キーフレーズも、用語 / フレーズ 108 とほぼ同様であるかどうかを判定する。同様である場合、このプロシージャは、ブロック 214 から継続され、このブロックにおいて、STS モジュール 112 は、この提案される用語リスト 110 をこのエンドユーザに送信する。そうでない場合には、このプロシージャは、ページ上のリファレンス「A」で示すように図 3 のブロック 302 から継続される。

【0046】

図 3 は、マルチセンスクエリについての関連語を提案するための例示のプロシージャ 300 を示している。プロシージャ 300 のオペレーションは、図 2 のプロシージャ 200 のオペレーションに基づいている。説明するために、このプロシージャのオペレーションでは、図 1 のコンポーネントに関連して説明している。(すべての参照番号は、そのコンポーネントが最初に導入される図面番号を用いて開始される)。ブロック 302 において、STS モジュール 112 は、用語クラスタ 138 から分類機構 140 を生成し、それはこの時点では高出現度数クエリ用語 120 に基づいている。ブロック 304 において、STS モジュール 112 は、低出現度数クエリ用語 122 を 1 つずつ検索エンジン 126 に送信し、対応する検索結果 130 を受信する。ブロック 306 において、STS モジュール 112 は、検索結果 130 から断片記述を抽出し (フィーチャ 134 を抽出しており)、そこから用語ベクトル 136 を生成する。ブロック 308 において、STS モジュール 112 は、トレーニングされた分類機構 140 に鑑みて低出現度数クエリ用語 122 から生成された用語ベクトル 136 を分類して、この低出現度数クエリ用語に基づいた各用語クラスタ 138 を生成する。

【0047】

ブロック 310 において、STS モジュール 112 は、用語 / フレーズ 108 とほぼ同

10

20

30

40

50

様であると判定されている低出現度数クエリ用語 1 2 2 に基づいた用語クラスタ 1 3 8 からのキーワード/キフレーズから、提案される用語リスト 1 1 0 を生成する。ブロック 3 1 2 において、S T S モジュール 1 1 2 は、この提案される用語リスト 1 1 0 をこのエンドユーザに送信する。

【 0 0 4 8 】

例示の動作環境

図 4 は、マルチセンスクエリについての関連語を提案するための図 1 のシステム 1 0 0 と図 2 および 3 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ティング環境 4 0 0 の一実施例を示している。例示のコンピューティング環境 4 0 0 は、適切なコンピューティング環境の一実施例にすぎず、本明細書中に説明されるシステムおよび方法の用途または機能の範囲についてのどのような限定をも示唆することを意図してはいない。また、コンピューティング環境 4 0 0 は、コンピューティング環境 4 0 0 に示すコンポーネントのうちのもの 1 つまたは組合せにも関連したどのような依存性または必要性を有するものとも解釈すべきではない。

【 0 0 4 9 】

本明細書中で説明してい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ティングシステム環境またはコンピューティングシステム構成を用いて動作することが可能である。使用するのに適したものとすることができるよく知られているコンピューティングシステム、コンピューティング環境、および/またはコンピューティング構成の実施例には、それだけには限定されないが、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、任意の上記システムまたはデバイスを含む分散コンピューティング環境などが含まれる。このフレームワークのコンパクトバージョンまたはサブセットバージョンはまた、ハンドヘルドコンピュータや他のコンピューティングデバイスなど限られたリソースのクライアント中で実装することもできる。本発明は、分散コンピューティング環境中でも実施され、ここでタスクは、通信ネットワークを介してリンクされるリモート処理デバイスによって実施される。分散コンピューティング環境においては、プログラムモジュールは、ローカルメモリストレージデバイスにも、リモートメモリストレージデバイスにも配置することができる。

【 0 0 5 0 】

図 4 を参照すると、マルチセンスクエリについての関連語を提案するための例示のシステムが、コンピュータ 4 1 0 の形態の汎用コンピューティングデバイスを含んでいる。コンピュータ 4 1 0 の以下に説明している態様は、P S S サーバ 1 0 2 (図 1) および/またはクライアントコンピューティングデバイス 1 0 6 の例示の実装形態である。コンピュータ 4 1 0 のコンポーネントは、それだけには限定されないが、1 つ (または複数) の処理装置 4 2 0、システムメモリ 4 3 0、およびこのシステムメモリを含めて様々なシステムコンポーネントを処理装置 4 2 0 に結合するシステムバス 4 2 1 を含むことができる。システムバス 4 2 1 は、メモリバスまたはメモリコントローラ、ペリフェラルバス、および様々なバスアーキテクチャのうちどれかを使用したローカルバスを含めて、いくつかのタイプのバス構造のうちどれにすることもできる。実施例として、限定するものではないが、かかるアーキテクチャには、I S A (I n d u s t r y S t a n d a r d A r c h i t e c t u r e 業界標準アーキテクチャ) バス、M C A (M i c r o C h a n n e l A r c h i t e c t u r e マイクロチャンネルアーキテクチャ) バス、E I S A (E n h a n c e d I S A 拡張 I S A) バス、V E S A (V i d e o E l e c t r o n i c s S t a n d a r d s A s s o c i a t i o n ビデオエレクトロニクス規格協会) ローカルバス、およびメザンバスとしても知られている P C I (P e r i p h e r a l C o m p o n e n t I n t e r c o n n e c t) バスが含まれる。

【 0 0 5 1 】

コンピュータ 4 1 0 は、一般的に様々なコンピュータ読取り可能媒体を含んでいる。コンピュータ読取り可能媒体は、コンピュータ 4 1 0 からアクセスすることができる使用可

10

20

30

40

50

能な任意の媒体とすることができ、揮発性媒体も不揮発性媒体も、着脱可能媒体も着脱不能媒体も共に含んでいる。実施例として限定するものではないが、コンピュータ読取り可能媒体は、コンピュータストレージ媒体および通信媒体を含むことができる。コンピュータストレージ媒体は、コンピュータ読取り可能命令、データ構造、プログラムモジュール、他のデータなどの情報のストレージのための任意の方法または技術で実装される揮発性媒体および不揮発性媒体、着脱可能媒体および着脱不能媒体を含んでいる。コンピュータストレージ媒体には、それだけには限定されないが、RAM、ROM、EEPROM、フラッシュメモリまたは他のメモリ技術、CD-ROM、DVD（デジタル多用途ディスク）また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または他の磁気ストレージデバイス、あるいは所望の情報を記憶するために使用することができ、コンピュータ410からアクセスすることができる他の任意の媒体が含まれる。

10

【0052】

通信媒体は、一般的に搬送波や他の搬送メカニズムなどの被変調データ信号の形のコンピュータ読取り可能命令、データ構造、プログラムモジュールまたは他のデータを実施し、任意の情報配信媒体を含んでいる。用語「被変調データ信号」は、信号に情報を符号化するようにその1つまたは複数の特性が設定または変更された信号を意味する。実施例として、限定するものではないが、通信媒体は、有線ネットワークや直接配線接続などの有線媒体と、音響、RF、赤外線、他の無線媒体などの無線媒体とを含んでいる。以上のうちの任意の組合せもまた、コンピュータ読取り可能媒体の範囲内に含まれるべきである。

20

【0053】

システムメモリ430は、ROM（読取り専用メモリ）431やRAM（ランダムアクセスメモリ）432など、揮発性および/または不揮発性のメモリの形態のコンピュータストレージ媒体を含んでいる。起動中などコンピュータ410内のエレメント間で情報を転送する助けをする基本ルーチンを含むBIOS（基本入出力システム）433は、一般的にROM431に記憶される。RAM432は、処理装置420にとって直接にアクセス可能な、または処理装置420によって現在動作させられている、あるいはその両方が行われるデータおよび/またはプログラムモジュールを一般的に含んでいる。実施例として、限定するものではないが、図4は、オペレーティングシステム434、アプリケーションプログラム435、他のプログラムモジュール436、およびプログラムデータ437を示している。一実装形態においては、コンピュータ410は、PSSサーバ102である。このシナリオにおいては、アプリケーションプログラム435は、検索用語提案モジュール112、および分類モジュール114を含んでいる。この同じシナリオにおいて、プログラムデータ437は、用語/フレーズ108、提案される用語リスト110、履歴クエリ116、検索クエリ128、検索結果130、複数の検索結果132、抽出済みのフィーチャ134、用語ベクトル136、キーワードクラスタ138、トレーニングされる分類機構140、および他のデータ124を含む。

30

【0054】

コンピュータ410は、他の着脱可能/着脱不能な、揮発性/不揮発性のコンピュータストレージ媒体を含むこともできる。実施例にすぎないが、図4は、着脱不能な不揮発性磁気媒体から情報を読み取りまたはそれに情報を書き込むハードディスクドライブ441、着脱可能な不揮発性磁気ディスク452から情報を読み取りまたはそれに情報を書き込む磁気ディスクドライブ451、およびCD-ROMや他の光媒体など着脱可能な不揮発性光ディスク456から情報を読み取りまたはそれに情報を書き込む光ディスクドライブ455を示している。例示の動作環境において使用することができる他の着脱可能/着脱不能な揮発性/不揮発性のコンピュータストレージ媒体には、それだけには限定されないが、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、ソリッドステートRAM、ソリッドステートROMなどが含まれる。ハードディスクドライブ441は、一般的にインターフェース440など着脱不能メモリイン

40

50

ターフェースを介してシステムバス421に接続され、磁気ディスクドライブ451および光ディスクドライブ455は、一般的にインターフェース450など着脱可能なメモリインターフェースによってシステムバス421に接続される。

【0055】

前述され図4に示されるこれらのドライブおよびその関連するコンピュータストレージ媒体は、コンピュータ410についてのコンピュータ読取り可能命令、データ構造、プログラムモジュールおよび他のデータの記憶域を提供する。図4において、例えばハードディスクドライブ441は、オペレーティングシステム444、アプリケーションプログラム445、他のプログラムモジュール446、およびプログラムデータ447を記憶するものとして示されている。これらのコンポーネントは、オペレーティングシステム434、アプリケーションプログラム435、他のプログラムモジュール436、およびプログラムデータ437と同じとすることもでき、また異なるものとするところにもできることに留意されたい。オペレーティングシステム444、アプリケーションプログラム445、他のプログラムモジュール446、およびプログラムデータ447には、これらが少なくとも異なるコピーであることを示すために、ここでは異なる番号が付与されている。

10

【0056】

ユーザは、キーボード462や、一般にマウス、トラックボールまたはタッチパッドと呼ばれるポインティングデバイス461などの入力デバイスを介してコンピュータ410にコマンドおよび情報を入力することができる。他の入力デバイス(図示せず)は、マイクロフォン、ジョイスティック、ゲームパッド、衛星パラボラアンテナ、スキャナなどを含むことができる。これらおよび他の入力デバイスは、このシステムバス421に結合されるユーザ入力インターフェース460を介して処理装置420にしばしば接続されるが、これは、パラレルポート、ゲームポート、USB(ユニバーサルシリアルバス)など他のインターフェースおよびバス構造によって接続することもできる。

20

【0057】

モニタ491または他のタイプのディスプレイデバイスもまた、ビデオインターフェース490などのインターフェースを介してシステムバス421に接続される。このモニタに追加して、コンピュータは、スピーカ497やプリンタ496などの他のペリフェラル出力デバイスを含むこともでき、これらは、出力ペリフェラルインターフェース495を介して接続することができる。

30

【0058】

コンピュータ410は、リモートコンピュータ480など1つまたは複数のリモートコンピュータに対する論理接続を使用してネットワーク環境中で動作する。リモートコンピュータ480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、ピアデバイス、または他の共通ネットワークノードとすることができ、またその特定の実装形態に応じて、コンピュータ410に関連した前述の要素の多くまたはすべてを含むこともできるが、メモリストレージデバイス481だけしか図4には示されていない。図4に示す論理接続は、LAN(ローカルエリアネットワーク)471およびWAN(ワイドエリアネットワーク)473を含んでいるが、他のネットワークを含むこともできる。かかるネットワーキング環境は、オフィス、企業規模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においては、一般的なものである。

40

【0059】

LANネットワーキング環境中で使用する際には、コンピュータ410は、ネットワークインターフェースまたはアダプタ470を介してLAN471に接続される。WANネットワーキング環境中で使用する際には、コンピュータ410は、一般的にインターネットなどのWAN473上で通信を確立するためのモデム472または他の手段を含んでいる。モデム472は、内蔵でも外付けでもよいが、ユーザ入力インターフェース460または他の適切なメカニズムを介してシステムバス421に接続することができる。ネットワーク環境においては、コンピュータ410に関連して示されるプログラムモジュール、またはその一部分は、リモートメモリストレージデバイスに記憶することができる。実施

50

例として、限定するものではないが、図4は、リモートアプリケーションプログラム485をメモリデバイス481上に存在するものとして示している。図に示すこのネットワーク接続は、例示的なものであり、コンピュータ間で通信リンクを確立する他の手段を使用することもできる。

【0060】

結論

マルチセンスクエリについての関連語を提案するためのシステムおよび方法について、構造的な特徴および/または方法的なオペレーションもしくは動作に特有の言語で説明してきたが、添付の特許請求の範囲中で定義される実装形態は、必ずしも説明したこの特定の特徴または動作だけに限定されずとは限らないことが理解されよう。したがって、この特定の特徴および動作は、この請求される内容を実装する例示の形態として開示されている。

10

【図面の簡単な説明】

【0061】

【図1】マルチセンスクエリについての関連語を提案するためのシステムの一例を示す図である。

【図2】マルチセンスクエリについての関連語を提案するためのプロシージャの一例を示す図である。

【図3】マルチセンスクエリについての関連語を提案するための、図2のオペレーションに基づいたオペレーションであるプロシージャの一例を示す図である。

20

【図4】マルチセンスクエリについての関連語を提案するための、本発明によるシステム、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コンピューティング環境の一例を示す図である。

【符号の説明】

【0062】

- 102 編集検証サーバ
- 106 クライアントコンピューティングデバイス
- 110 提案される用語リスト
- 112 検索用語提案モジュール
- 114 分類モジュール
- 116 履歴クエリ
- 120 高FOOクエリ用語
- 122 低FOOクエリ用語
- 124 他のデータ
- 132 検索エンジン
- 128 1つ(または複数)のクエリログ
- 132 検索結果
- 134 抽出されたフィーチャ
- 136 1つ(または複数)の用語ベクトル
- 138 用語クラスタ
- 140 トレーニングされた分類機構
- 420 1つ(または複数)の処理装置
- 421 システムバス
- 430 システムメモリ
- 434 オペレーティングシステム
- 435 アプリケーションプログラム
- 436 他のプログラムモジュール
- 437 プログラムデータ
- 440 着脱不能な揮発性メモリインターフェース
- 444 オペレーティングシステム

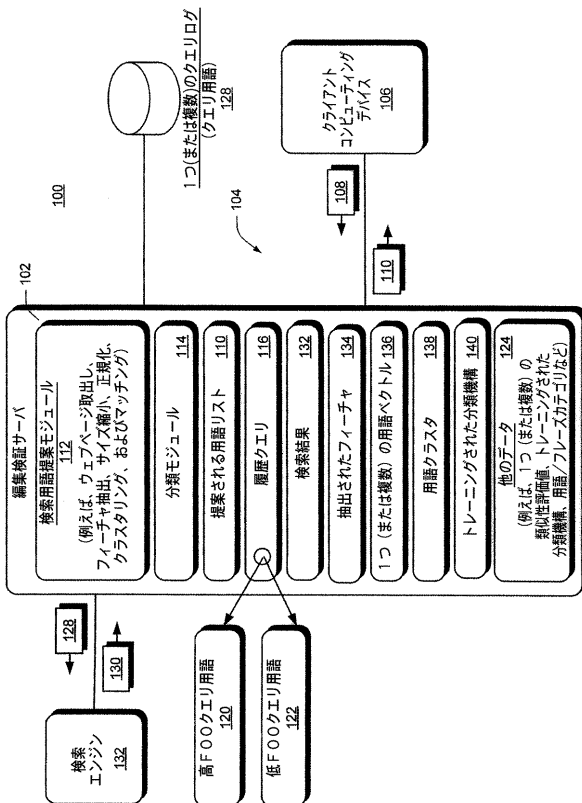
30

40

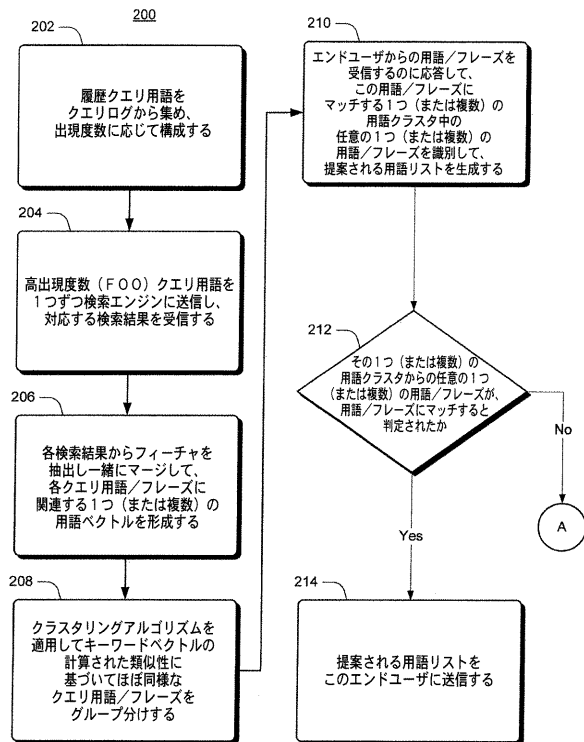
50

- 4 4 5 アプリケーションプログラム
- 4 4 6 他のプログラムモジュール
- 4 4 7 プログラムデータ
- 4 5 0 着脱可能な不揮発性メモリインターフェース
- 4 6 0 ユーザ入力インターフェース
- 4 6 1 マウス
- 4 6 2 キーボード
- 4 7 0 ネットワークインターフェース
- 4 7 1 ローカルエリアネットワーク
- 4 7 2 モデム
- 4 7 3 ワイドエリアネットワーク
- 4 8 0 リモートコンピュータ
- 4 8 5 リモートアプリケーションプログラム
- 4 9 0 ビデオインターフェース
- 4 9 1 モニタ
- 4 9 4 入力ペリフェラルインターフェース
- 4 9 5 出力ペリフェラルインターフェース
- 4 9 6 プリンタ
- 4 9 7 スピーカ

【図 1】



【図 2】



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 鈴木 和樹

- (56)参考文献 特開2003 - 233684 (JP, A)
特表2002 - 518748 (JP, A)
Vijay V. Raghavan、外1名、ON THE REUSE OF PAST OPTIMAL QUERIES, Proceedings of the 18
th annual international ACM SIGIR conference on Research and development in informatio
n retrieval(SIGIR'95)[online], 1995年, p. 344 - 350, [DL from ACM Digital Li
brary]
Myoung-Cheol Kim、外1名、A comparison of collocation-based similarity measures in que
ry expansion, Information Processing and Management[online], 1999年, 第35巻, 第
1号, p. 19 - 30, [DL from Science Direct]
Yonggang Qiu、外1名、Concept Based Query Expansion, Proceedings of the 16th annual in
ternational ACM SIGIR conference on Research and development in information retrieval(
SIGIR'93)[online], 1993年, p. 160 - 169, [DL from ACM Digital Library]

(58)調査した分野(Int.Cl., DB名)

G06F 17/30

THE ACM DIGITAL LIBRARY

Science Direct