

(19) 日本国特許庁 (JP)

(12) 特 許 公 報 (B2)

(11) 特許番号

特許第6284643号  
(P6284643)

(45) 発行日 平成30年2月28日 (2018. 2. 28)

(24) 登録日 平成30年2月9日 (2018. 2. 9)

(51) Int. Cl.

F I

G 0 6 F 17/30 (2006.01)

G 0 6 F 17/30 2 1 0 A

G 0 6 F 17/30 1 6 0 F

請求項の数 20 (全 20 頁)

|               |                               |           |                     |
|---------------|-------------------------------|-----------|---------------------|
| (21) 出願番号     | 特願2016-536850 (P2016-536850)  | (73) 特許権者 | 516164586           |
| (86) (22) 出願日 | 平成26年12月1日 (2014. 12. 1)      |           | キューベース リミテッド ライアビリテ |
| (65) 公表番号     | 特表2016-541069 (P2016-541069A) |           | ィ カンパニー             |
| (43) 公表日      | 平成28年12月28日 (2016. 12. 28)    |           | アメリカ合衆国 ヴァージニア州 201 |
| (86) 国際出願番号   | PCT/US2014/067918             |           | 91 レストン サンライズ ヴァレー  |
| (87) 国際公開番号   | W02015/084724                 |           | ドライヴ 12018 スイート 300 |
| (87) 国際公開日    | 平成27年6月11日 (2015. 6. 11)      | (74) 代理人  | 100086771           |
| 審査請求日         | 平成29年11月22日 (2017. 11. 22)    |           | 弁理士 西島 孝喜           |
| (31) 優先権主張番号  | 61/910, 739                   | (74) 代理人  | 100088694           |
| (32) 優先日      | 平成25年12月2日 (2013. 12. 2)      |           | 弁理士 弟子丸 健           |
| (33) 優先権主張国   | 米国 (US)                       | (74) 代理人  | 100094569           |
| 早期審査対象出願      |                               |           | 弁理士 田中 伸一郎          |
|               |                               | (74) 代理人  | 100067013           |
|               |                               |           | 弁理士 大塚 文昭           |

最終頁に続く

(54) 【発明の名称】 非構造化テキストにおける特徴の曖昧性除去方法

(57) 【特許請求の範囲】

【請求項 1】

サーチ質問をエンドユーザ装置から受信することに応答して、

システムのノードにより、1つ以上の抽出された特徴に一致する1つ以上の候補レコードを識別するために共起特徴を含む候補レコードのセットをサーチし、候補レコードに一致する抽出された特徴は、一次特徴であり、前記ノードはインメモリデータベースをホストするメインメモリを含み、前記インメモリデータベースはクラスターの知識ベースを格納し、各クラスターは独特の識別子（独特のID）を伴う曖昧性除去された一次特徴及び関連付けられた二次特徴を含み、

ノードにより、抽出された特徴の各々を1つ以上のマシン発生トピック識別子（トピックID）と関連付け、

ノードにより、トピックIDの関連度に基づき一次特徴の各々を互いに曖昧性除去し、ノードにより、トピックIDの関連度に基づき各一次特徴に関連した二次特徴のセットを識別し、

ノードにより、トピックIDの関連度に基づき二次特徴の関連セットにおける二次特徴の各々から一次特徴の各々を曖昧性除去し、

前記インメモリデータベースから前記知識ベースからデータが検索されるときに、リアルタイムで、ノードにより、各一次特徴を二次特徴の関連セットにリンクして、新たなクラスターを形成し、

ノードのインメモリデータベースの曖昧性除去モジュールにより、曖昧性除去された一

10

20

次特徴を伴う既存の知識クラスターへの比較的一致するスコアの指定により前記新たなクラスターの各々が既存の知識ベースクラスターに一致するかどうか決定し、

一致があるときには、既存の知識ベースクラスターにおける各一致する一次特徴に対応する既存の独特のIDを決定しそして前記新たなクラスターを含むように既存の知識ベースクラスターを更新し、及び

一致がないときには、新たな知識ベースクラスターを生成し、そしてその新たな知識ベースクラスターの一次特徴に新たな独特のIDを指定し、及び

既存の独特のID及び新たな独特のIDの一方を一次特徴として前記エンドユーザ装置へ送出する、  
ことを含む方法。

10

【請求項2】

ノードにより、抽出された特徴に一致する候補レコードの各々を比較し、及びノードにより、その比較に基づいて前記抽出された特徴の各々に重み付けされた一致スコア結果を指定する、ことを更に含む、請求項1に記載の方法。

【請求項3】

ノードにより、抽出された特徴の各々を、重み付けされた特徴属性のセットに関連付けることを更に含む、請求項2に記載の方法。

【請求項4】

ノードにより、1つ以上の重み付けされた特徴属性に基づいて抽出された特徴の各々の関連度を決定することを更に含む、請求項3に記載の方法。

20

【請求項5】

1つ以上の抽出された特徴をノードの抽出モジュールにより確認及び抽出し、1つ以上の抽出された特徴において1つ以上の一次特徴を識別し、及び

ノードの抽出モジュールにより、抽出された特徴の各々をデータベースに記憶する、ことを更に含む、請求項1に記載の方法。

【請求項6】

ノードの抽出モジュールにより、各特徴に抽出確度スコアを指定することを更に含む、請求項5に記載の方法。

【請求項7】

各々の一次特徴は、1つ以上の特徴属性のセットに関連付けられる、請求項1に記載の方法。

30

【請求項8】

特徴属性は、トピックID、ドキュメント識別子(ドキュメントID)、特徴タイプ、特徴名、信頼性スコア、及び特徴位置より成るグループから選択される、請求項7に記載の方法。

【請求項9】

各関連特徴は、予め定義されたクラスターハイアラキーに従って下位順序特徴のセットに関連付けられる、請求項1に記載の方法。

【請求項10】

ノードにより、候補レコードのセットの曖昧キーサーチを遂行することを更に含む、請求項1に記載の方法。

40

【請求項11】

ノードのリンクオンザフライモジュールにより、関連トピックIDの共起及び1つ以上の特徴属性に基づいて2つ以上のデータソースをリンクすることを更に含む、請求項7に記載の方法。

【請求項12】

ノードにより、データソースにおける抽出された特徴が第2データソースにおいて共起するかどうかを、その抽出された特徴を第2データソースにおける特徴と比較することで決定し、及び

ノードにより、前期比較に基づいてデータソースの各々をリンクする、

50

ことを更に含む、請求項 1 に記載の方法。

【請求項 1 3】

ノードにより、異なるデータソースからの抽出された特徴の共起を分析して、抽出された特徴の曖昧性除去の精度を改善することを更に含む、請求項 1 に記載の方法。

【請求項 1 4】

ノードにより、1 つ以上の新たなデータソースを連続的に受け取り、  
ノードにより、1 つ以上の抽出される特徴を連続的に抽出し、  
ノードにより、1 つ以上の抽出された特徴において候補サーチを連続的に遂行し、  
ノードにより、抽出された特徴を連続的に曖昧性除去し、及び  
ノードにより、抽出された特徴を 1 つ以上の新たなクラスターへ連続的にリンクする、  
ことを更に含む、請求項 1 に記載の方法。

【請求項 1 5】

コンピュータ実行可能なインストラクションが記憶された非一時的コンピュータ読み取り可能な媒体であって、プロセッサによって実行されると、

サーチ質問をエンドユーザ装置から受信することに応答して、

システムのノードにより、1 つ以上の抽出された特徴に一致する 1 つ以上の候補レコードを識別するために共起特徴を含む候補レコードのセットをサーチし、前記ノードはインメモリデータベースをホストするメインメモリを含み、前記インメモリデータベースはクラスターの知識ベースを格納し、各クラスターは独特の識別子（独特の ID）を伴う曖昧性除去された一次特徴及び関連付けられた二次特徴を含み、

ノードにより、抽出された特徴の各々を 1 つ以上のマシン発生トピック識別子（トピック ID）と関連付け、

ノードにより、トピック ID の関連度に基づき一次特徴の各々を互いに曖昧性除去し、

ノードにより、トピック ID の関連度に基づき各一次特徴に関連した二次特徴のセットを識別し、

ノードにより、トピック ID の関連度に基づき二次特徴の関連セットにおける二次特徴の各々から一次特徴の各々を曖昧性除去し、

前記インメモリデータベースから前記知識ベースからデータが検索されるときに、リアルタイムで、ノードにより、各一次特徴を二次特徴の関連セットにリンクして、新たなクラスターを形成し、

ノードのインメモリデータベースの曖昧性除去モジュールにより、曖昧性除去された一次特徴を伴う既存の知識クラスターへの比較的一致するスコアの指定により前記新たなクラスターの各々が既存の知識ベースクラスターに一致するかどうか決定し、

一致があるときには、既存の知識ベースクラスターにおける各一致する一次特徴に対応する既存の独特の ID を決定しそして前記新たなクラスターを含むように既存の知識ベースクラスターを更新し、及び

一致がないときには、新たな知識ベースクラスターを生成し、そしてその新たな知識ベースクラスターの一次特徴に新たな独特の ID を指定し、及び

既存の独特の ID 及び新たな独特の ID の一方を一次特徴として前記エンドユーザ装置へ送出する、

ことを含む機能が実行される、コンピュータ実行可能なインストラクションが記憶された非一時的コンピュータ読み取り可能な媒体。

【請求項 1 6】

前記インストラクションは、更に、ノードにより、抽出された特徴に一致する候補レコードの各々を比較し、及びノードにより、その比較に基づいて前記抽出された特徴の各々に重み付けされた一致スコア結果を指定する、ことを含む、請求項 1 5 に記載の非一時的コンピュータ読み取り可能な媒体。

【請求項 1 7】

前記インストラクションは、更に、ノードにより、抽出された特徴の各々を、重み付けされた特徴属性のセットに関連付けることを含む、請求項 1 6 に記載の非一時的コンピュ

10

20

30

40

50

ータ読み取り可能な媒体。

【請求項 18】

前記インスタレーションは、更に、ノードにより、1つ以上の重み付けされた特徴属性に基づいて抽出された特徴の各々の関連度を決定することを含む、請求項 17 に記載の非一時的コンピュータ読み取り可能な媒体。

【請求項 19】

前記インスタレーションは、更に、

ノードの抽出モジュールにより、1つ以上の抽出された特徴を確認し及び抽出し、その1つ以上の抽出された特徴において1つ以上の一次特徴を識別し、及び

ノードの抽出モジュールにより、抽出された特徴の各々をデータベースに記憶する、  
ことを含む、請求項 15 に記載の非一時的コンピュータ読み取り可能な媒体。

10

【請求項 20】

前記インスタレーションは、更に、ノードの抽出モジュールにより、各特徴に抽出確度スコアを指定することを含む、請求項 19 に記載の非一時的コンピュータ読み取り可能な媒体。

【発明の詳細な説明】

【技術分野】

【0001】

本発明は、一般的に、データマネージメントに関するもので、より詳細には、ネットワークを経て受け取ったソースアイテムから資料を抽出しそして記憶するデータマネージメントシステム及び方法に関する。

20

【背景技術】

【0002】

ネットワークのようなソースを含む大きなドキュメント集合体におけるエンティティ（即ち、人々、位置、組織）に関する情報のサーチは、しばしば曖昧なものであり、不正確なテキスト処理機能、知識抽出中の特徴の不正確な関連付け、ひいては、不正確なデータ分析を招くことがある。

【0003】

最新のシステムは、Page Rank 及びハイパーリンク誘起トピックサーチ（HITS）アルゴリズムのような多数のアルゴリズムにおけるリンケージベースのクラスター化及びランク付けを使用している。この解決策及び関連解決策の背景にある基本的アイデアは、既存のリンクが、典型的に、関連ページ又はコンセプト間に存在することである。クラスター化ベースの技術の制約は、エンティティを曖昧性除去するのに必要なコンテキスト情報がコンテキストに存在しないことが時々あって、不正確な曖昧性除去結果を招くことである。同様に、同じ又は表面的に同様のコンテキストにおける異なるエンティティに関するドキュメントが不正確にクラスター化されることもある。

30

【0004】

他のシステムは、エンティティの1つ以上の外部ディクショナリ（又は知識ベース）を参照することによりエンティティの曖昧性除去を試みる。そのようなシステムでは、エンティティのコンテキストがディクショナリ内の考えられる一致エンティティと比較され、そして最も近い一致が返送される。現在のディクショナリベース技術に関連した制約は、いつでもエンティティの数が増加し、それ故、世界中の全てのエンティティの表現を含むディクショナリがないことから生じる。従って、ドキュメントのコンテキストがディクショナリのエンティティに一致する場合に、この技術では、ディクショナリ内の最も類似したエンティティだけが識別され、ディクショナリ以外にある正しいエンティティが必ずしも識別されない。

40

【発明の概要】

【発明が解決しようとする課題】

【0005】

ほとんどの方法は、曖昧性除去プロセスにおいてエンティティ及びキーフレーズだけを

50

使用する。それ故、正確なデータ分析を行うことのできる正確なエンティティ曖昧性除去技術が依然要望されている。

【課題を解決するための手段】

【0006】

ある実施形態では、特徴を曖昧性除去する方法について述べる。この方法は、1つ以上の特徴抽出モジュール、1つ以上の曖昧性除去モジュール、1つ以上のスコア付けモジュール、及び1つ以上のリンクングモジュールのような複数のモジュールを含む。

【0007】

特徴の曖昧性除去は、特徴の周囲ドキュメントからトピックを抽出し、レイテントディリクレアロケーション(MC-LDA)トピックモデルのマルチコンポーネント拡張を使用することにより、一部分、サポートされる。ここで、各コンポーネントは、既存の知識ベースに記憶されるか又は到来するドキュメントにおいて抽出される各二次的特徴に関してモデリングされる。更に、リンクング又は曖昧性除去プロセスは、MC-LDAからのトピック推論としてモデリングされ、これは、MC-LDAトレーニングの間に自動重み推定を与え、そして推論中にそれを容易に適用する。

【0008】

この規範的方法は、ドキュメントリンクングを考慮せずに達成できるものを越えてエンティティ曖昧性除去の精度を改善する。ドキュメントリンケージを考慮し、リンクにより暗示されるドキュメント及びエンティティ関係を考えることで良好な曖昧性除去を行うことができる。

【0009】

ある実施形態において、方法は、インメモリデータベースをホストするシステムのノードにより、1つ以上の抽出された特徴に一致する1つ以上の候補を識別するために候補レコードのセットをサーチし、候補に一致する抽出された特徴は、一次特徴であり；ノードにより、抽出された特徴の各々を1つ以上のマシン発生トピック識別子(トピックID)と関連付け；ノードにより、トピックIDの関連度に基づき一次特徴の各々を互いに曖昧性除去し；ノードにより、トピックIDの関連度に基づき各一次特徴に関連した二次特徴のセットを識別し；ノードにより、トピックIDの関連度に基づき二次特徴の関連セットにおける二次特徴の各々から一次特徴の各々を曖昧性除去し；ノードにより、各一次特徴を二次特徴の関連セットにリンクして新たなクラスターを形成し；ノードにより、新たなクラスターが既存の知識ベースクラスターに一致するかどうか決定し、一致があるときには、インメモリデータベースサーバーコンピュータの曖昧性除去モジュールにより、知識ベースクラスターにおける各一致する一次特徴に対応する既存の独特の識別子(独特のID)を決定しそしてその新たなクラスターを含むように知識ベースクラスターを更新し、及び一致がないときには、ノードにより、新たな知識ベースクラスターを生成し、そしてその新たな知識ベースクラスターの一次特徴に新たな独特のIDを指定し；及びノードにより、既存の独特のID及び新たな独特のIDの一方を一次特徴として送出する；ことを含む。

【0010】

別の実施形態において、非一時的コンピュータ読み取り可能な媒体に記憶されるコンピュータ実行可能なインストラクションは、インメモリデータベースをホストするシステムのノードにより、1つ以上の抽出された特徴に一致する1つ以上の候補を識別するために候補レコードのセットをサーチし、候補に一致する抽出された特徴は、一次特徴であり；ノードにより、抽出された特徴の各々を1つ以上のマシン発生トピック識別子(トピックID)と関連付け；ノードにより、トピックIDの関連度に基づき一次特徴の各々を互いに曖昧性除去し；ノードにより、トピックIDの関連度に基づき各一次特徴に関連した二次特徴のセットを識別し；ノードにより、トピックIDの関連度に基づき二次特徴の関連セットにおける二次特徴の各々から一次特徴の各々を曖昧性除去し；ノードにより、各一次特徴を二次特徴の関連セットにリンクして新たなクラスターを形成し；ノードにより、新たなクラスターが既存の知識ベースクラスターに一致するかどうか決定し、一致がある

ときには、ノードにより、知識ベースクラスターにおける各一致する一次特徴に対応する既存の独特の識別子（独特のID）を決定しそしてその新たなクラスターを含むように知識ベースクラスターを更新し、及び一致がないときには、新たな知識ベースクラスターを生成し、そしてその新たな知識ベースクラスターの一次特徴に新たな独特のIDを指定し；及びノードにより、既存の独特のID及び新たな独特のIDの一方を一次特徴として送出する；ことを含む。

【0011】

一実施形態の付加的な特徴及び効果は以下の説明に述べられ、そして一部分はその説明から明らかとなる。本発明の目的及び他の効果は、以下の説明の規範的实施形態、特許請求の範囲及び添付図面に特に指摘された構造により実現され且つ達成される。

10

【0012】

以上の一般的な説明及び以下の詳細な説明は、どちらも、規範的な説明に過ぎず、請求の範囲に述べる本発明の更なる説明を与えるものであることを理解されたい。

【0013】

本開示は、添付図面を参照することにより良く理解することができる。添付図面は、本明細書の一部を構成するもので、本発明の実施形態を示し、そして明細書と共に本発明を説明するものである。図面中のコンポーネントは、必ずしも、正しい縮尺ではなく、むしろ、本開示の原理を例示する際に強調されている。図中、参照番号は、異なる図面全体を通して対応部分を示している。

【図面の簡単な説明】

20

【0014】

【図1】規範的实施形態により、非構造化テキストにおける特徴を曖昧性除去する方法のフローチャートである。

【図2】規範的实施形態により、特徴を曖昧性除去する方法に使用される曖昧性除去モジュールにより遂行されるステップのフローチャートである。

【図3】規範的实施形態により、特徴を曖昧性除去する方法に使用されるリンクオンザフライモジュールにより遂行されるステップのフローチャートである。

【図4】規範的实施形態により、特徴を曖昧性除去する方法の実施に使用されるシステムを示す図である。

【図5】規範的实施形態によるマルチコンポーネント条件独立のレイテンシディリクレアロケーション（MC-LDA）トピックモデルのグラフ表示である。

30

【図6】規範的实施形態によるマルチコンポーネント条件独立のレイテンシディリクレアロケーション（MC-LDA）トピックモデルのギブスサンプリング方程式の実施形態を示す。

【図7】規範的实施形態によるマルチコンポーネント条件独立のレイテンシディリクレアロケーション（MC-LDA）トピックモデルにおけるトレーニング及び推論のための確率変化推論アルゴリズムの実施形態を示す。

【図8】規範的实施形態によるマルチコンポーネント条件独立のレイテンシディリクレアロケーション（MC-LDA）トピックモデルのためのサンプルトピックを示すテーブルである。

40

【発明を実施するための形態】

【0015】

定義

ここで使用する次の用語は、次のような定義を有する。

【0016】

「ドキュメント」は、出発点及び終了点を有する情報の個別の電子的表現を指す。

【0017】

「マルチドキュメント」は、トークン、異なる形式の名前付きエンティティ、及び個別の「bag-of-surface-forms」コンポーネントに編成されるキーフレーズを伴うドキュメントを指す。

50

## 【 0 0 1 8 】

「データベース」は、1つ以上の集合体を記憶するのに適し且つ1つ以上の質問を処理するのに適したクラスター及びモジュールの組み合わせを含むシステムを指す。

## 【 0 0 1 9 】

「コーパス」は、1つ以上のドキュメントの集合体を指す。

## 【 0 0 2 0 】

「生のコーパス」又は「ドキュメントストリーム」は、新たなドキュメントがネットワークにアップロードされるときに恒常的に供給されるコーパスを指す。

## 【 0 0 2 1 】

「特徴(Features)」は、ドキュメントから少なくとも一部分導出される情報である。

10

## 【 0 0 2 2 】

「特徴属性」は、特徴に関連したメタデータ、例えば、とりわけ、ドキュメントにおける特徴の位置、信頼スコアを指す。

## 【 0 0 2 3 】

「クラスター」は、特徴の集合体を指す。

## 【 0 0 2 4 】

「エンティティ知識ベース」は、特徴/エンティティを含むベースを指す。

## 【 0 0 2 5 】

「リンクオンザフライモジュール」又は「リンクO T F」は、生のコーパスが更新されるにつれてデータを更新するリンキングモジュールを指す。

20

## 【 0 0 2 6 】

「メモリ」は、十分に高い速度で情報を記憶し且つその情報を検索するのに適したハードウェアコンポーネントを指す。

## 【 0 0 2 7 】

「モジュール」は、1つ以上の定義されたタスクを実行するのに適したコンピュータソフトウェアコンポーネントを指す。

## 【 0 0 2 8 】

「センチメント(Sentiment)」は、ドキュメント、ドキュメントの一部分、又は特徴に関連した客観的評価を指す。

## 【 0 0 2 9 】

30

「トピック」は、コーパスから少なくとも一部分導出されるセマティック情報のセットを指す。

## 【 0 0 3 0 】

「トピック識別子」又は「トピックID」は、トピックの特定インスタンスを指す識別子である。

## 【 0 0 3 1 】

「トピック集合体」は、コーパスから導出されるトピックの特定セットを指し、各トピックは、独特の識別子(独特のID)を有する。

## 【 0 0 3 2 】

「トピック分類」は、特定のトピック識別子をドキュメントの特徴として指定することを指す。

40

## 【 0 0 3 3 】

「質問」は、1つ以上の適当なデータベースから情報を検索するための要求を指す。

## 【 0 0 3 4 】

詳細な説明

添付図面に各々示された好ましい実施形態を以下に詳細に説明する。上述した実施形態は、例示に過ぎない。当業者であれば、ここに述べる特定の実施例について、本発明の範囲内で、多数の別のコンポーネント及び実施形態に置き換えることが認識されよう。

## 【 0 0 3 5 】

本開示は、非構造化テキストにおける特徴を曖昧性除去する方法について述べる。規範

50

的な実施形態は、本開示に従って特徴を曖昧性除去する慣習について述べるが、ここに述べるシステム及び方法は、本開示の範囲内で適当に使用するように構成できることが意図される。

【0036】

既存の知識ベースは、曖昧さのない特徴及びそれに関連した特徴を含み、これは、信頼性の低いテキスト分析を招く。本開示の観点では、特徴及びエンティティの曖昧性除去精度を高め、それ故、テキスト分析の精度を高める。

【0037】

一実施形態によれば、特徴を曖昧性除去するここに開示する方法は、初期データコーパスに使用されて、ドキュメント取り込み及び特徴抽出を遂行し、初期コーパスに含まれた各ドキュメントに対してトピック分類及び他のテキスト分析を行えるようにする。各特徴は、とりわけ、ドキュメントの名前、タイプ、位置情報、及び信頼性スコアとして識別されそして記録される。

10

【0038】

図1は、非構造化テキストにおける特徴を曖昧性除去する複数のステップを示す方法100のフローチャートである。一実施形態によれば、特徴曖昧性除去方法100は、既存の知識ベースにおいて新たなドキュメント入力ステップ102が実行されるときに開始する。ドキュメントに対して特徴抽出ステップ104が遂行される。一実施形態によれば、特徴は、とりわけ、トピック識別子(トピックID)、ドキュメント識別子(ドキュメントID)、特徴のタイプ、特徴の名前、信頼性スコア及び特徴の位置のような異なる特徴属性に関連している。

20

【0039】

種々の実施形態によれば、ステップ102のドキュメント入力は、大量コーパス又は生のコーパス(インターネット又はネットワーク接続のコーパスのような)から供給され、これは、次いで、毎秒供給される。

【0040】

異なる実施形態によれば、特徴抽出ステップ104の間に、ドキュメント入力ステップ102の非構造化テキストを分析するために1つ以上の特徴確認及び抽出アルゴリズムが使用される。抽出された各特徴にスコアが指定される。そのスコアは、正しい属性で正しく抽出される特徴の確度レベルを指示する。

30

【0041】

加えて、特徴抽出ステップ104の間に、ステップ102のドキュメント入力から1つ以上の一次特徴が識別される。各一次特徴は、特徴属性のセット及び1つ以上の二次特徴に関連付けられる。各二次特徴は、特徴属性のセットに関連付けられる。ある実施形態では、1つ以上の二次特徴は、特徴属性のそれ自身のセットを各々有する1つ以上の三次特徴を有する。

【0042】

特徴属性を考慮して、ステップ102のドキュメント入力内の各特徴の相対的重み又は関連度が決定される。加えて、重み付けされたスコア付けモデルを使用して、特徴と特徴との間の関連付けの関連度が決定される。

40

【0043】

特徴抽出ステップ104に続いて、ステップ102のドキュメント入力から抽出された特徴及びそれに関連した全ての情報は、ステップ106においてMemDBに特徴を含ませる間に、特徴曖昧性除去要求ステップ108の一部として、インメモリデータベース(MemDB)にロードされる。

【0044】

一実施形態では、MemDBは、図1から8に関連して述べるステップを実行する1つ以上のプロセッサを有する曖昧性除去コンピュータサーバー環境の一部を形成する。ある実施形態では、MemDBは、1つ以上のサーチコントローラ、複数のサーチノード、圧縮データの集合体、及び曖昧性除去サブモジュールを含むコンピュータモジュールであ

50



る。１つのサーチコントローラが１つ以上のサーチノードに選択的に関連付けられる。各サーチノードは、圧縮データの集合体を通して曖昧キーワードを独立して遂行し、そしてスコア付けされた結果のセットをそれに関連したサーチコントローラへ返送することができる。

#### 【 0 0 4 5 】

特徴の曖昧性除去ステップ 1 0 8 は、M e m D B 内の曖昧性除去サブモジュールにより遂行される。特徴の曖昧性除去 1 0 8 プロセスは、マシンで発生されるトピック I D を含み、これは、特徴、ドキュメント、又はコーパスを分類するのに使用される。個々の特徴及び特定のトピック I D の関連度は、曖昧性除去アルゴリズムを使用して決定される。あるドキュメントにおいて、そのドキュメント内の特徴の異なる発生のコンテキストに基づき、１つ以上のトピック I D に同じ特徴が関連付けられる。

10

#### 【 0 0 4 6 】

あるドキュメントから抽出された特徴（同じトピック、接近用語及びエンティティ、キーワード、イベント及びファクト）のセットは、異なるドキュメントにわたる２つ以上の特徴が単一の特徴である場合、又はそれらが別々の特徴である場合に、ある精度レベルで定義する曖昧性除去アルゴリズムを使用して、他のドキュメントからの特徴のセットと比較される。ある例では、データベースにおけるドキュメントの集合体にわたる２つ以上の特徴の共起を分析して、特徴曖昧性除去プロセス 1 0 8 の精度を改善する。ある実施形態では、全体的スコア付けアルゴリズムを使用して、特徴が同じである確率を決定する。

#### 【 0 0 4 7 】

20

ある実施形態では、特徴曖昧性除去プロセス 1 0 8 の一部分として、M e m D B 内に知識ベースが発生される。この知識ベースは、関連する曖昧性除去された一次特徴及びそれに関連する二次特徴のクラスターを一時的に記憶するのに使用される。新たなドキュメントが M e m D B にロードされたときに、曖昧性除去された新たな特徴セットを既存の知識ベースと比較し、特徴と特徴との関係を決定し、そして新たな特徴と既に抽出された特徴との間に一致があるかどうか決定する。

#### 【 0 0 4 8 】

比較された特徴が一致する場合には、知識ベースが更新され、一致する特徴の特徴 I D がユーザ及び／又は要求側アプリケーション又はプロセスへ返送され、そして更に、一致の頻度に基づいて目立った手段を特徴 I D と共に取り付けることができ、これは、所与のコーパスにおいてその人気指数を捕らえるものである。比較された特徴が既に抽出された特徴のいずれとも一致しない場合には、曖昧性除去されたエンティティ又は特徴に独特の特徴 I D が指定され、その独特の特徴 I D は、特徴を定義するクラスターに関連付けられそして M e m D B の知識ベース内に記憶される。その後、ステップ 1 1 0 において、曖昧性除去された特徴の特徴 I D がシステムインターフェイスを通してソースへ返送される。ある実施形態では、曖昧性除去された特徴の特徴 I D は、二次特徴、特徴のクラスター、関連特徴属性、又は他の要求データを含む。特徴曖昧性除去ステップ 1 0 8 に対して使用される曖昧性除去サブモジュールを、図 2 について以下に詳細に述べる。

30

#### 【 0 0 4 9 】

##### 曖昧性除去サブモジュール

40

図 2 は、一実施形態により、方法 1 0 0（図 1）の特徴曖昧性除去ステップ 1 0 8 の非構造化テキストに使用される曖昧性除去サブモジュールにより遂行されるプロセス 2 0 0 のフローチャートである。曖昧性除去プロセス 2 0 0 は、図 1 のステップ 1 0 6 において M e m D B に特徴を含ませた後に始まる。ステップ 2 0 2 において与えられる抽出された特徴は、ステップ 2 0 4 において候補サーチを遂行するのに使用され、抽出された特徴についてのサーチは、共起特徴を含めて全ての候補レコードを通して遂行される。

#### 【 0 0 5 0 】

種々の実施形態によれば、候補は、特徴の曖昧性除去プロセス 1 0 8 に使用される関連二次特徴のセットを伴う一次特徴である。

#### 【 0 0 5 1 】

50

曖昧性除去結果は、トピックIDの共起とトピックIDの中の関連度とにより改善される。トピックIDの関連度は、異なるトピックモデルにわたるものであっても、トピックIDが指定された大きなコーパスから発見することができる。関連トピックIDをレコードリンケージステップ206の間に使用して、厳密なトピックIDを含まないが1つ以上の関連トピックIDを含むドキュメントへのリンケージを与えることができる。この解決策は、レコードリンケージステップ206に含まれるべき関連特徴のリコールを改善し、そしてあるケースでは、曖昧性除去結果を改善する。

#### 【0052】

潜在的に関連するドキュメントのセットが識別され、そしてそれらのドキュメント内の関連する一次及び二次特徴が抽出されると、特徴の属性、同じドキュメント（意義のあるコンテキスト）の特徴と特徴との間の関係、特徴の相対的重み、及び他の変数をレコードリンケージプロセス206の間に使用して、それらのドキュメントにわたる一次及び二次特徴を曖昧性除去する。次いで、各レコードを他のレコードにリンクして、曖昧性除去された一次特徴及びそれらの関連する二次特徴のクラスターを決定する。レコードリンケージ206に使用されるアルゴリズムは、マイニング非構造化データセットのスペルエラー又は翻字及び他の課題を克服することができる。

#### 【0053】

クラスター比較ステップ208は、比較的一致するスコアを、曖昧性除去された特徴のクラスターに指定することを含み、異なるアプリケーションに対して異なる受け入れスレッショールドが定義される。定義された精度レベルは、どのスコアが肯定的一致サーチと考えられそしてどのスコアが否定的一致サーチと考えられるか決定する（ステップ210）。各新たなクラスターは、独特のIDが与えられ、そして知識ベースに一時的に記憶される。各新たなクラスターは、曖昧性除去された新たな一次特徴及び二次特徴のセットを含む。新たなクラスターが、知識ベースに既に記憶されているクラスターに一致する場合には、システムは、ステップ212において知識ベースを更新し、そしてユーザ及び/又は要求側アプリケーション又はプロセスへの一致特徴IDの返送がステップ214において遂行される。知識ベースの更新212は、1つの一次特徴への付加的な二次特徴の関連付け、或いは一次又は二次特徴に以前に関連付けされていない特徴属性の追加を意味する。

#### 【0054】

評価されているクラスターに、肯定的一致サーチ210のスレッショールドより低いスコアが指定された場合には、システムは、ステップ216において、クラスターの一次特徴に独特のID指定を遂行し、そしてステップ212において、知識ベースを更新する。その後、システムは、一致ID返送プロセス214を遂行する。レコードリンケージステップ206は、図3を参照して更に詳細に説明する。

#### 【0055】

##### リンクオンザフライサブモジュール

図3は、一実施形態により、特徴を曖昧性除去する方法100に使用されるリンクオンザフライ（リンクOTF）サブモジュールにより遂行されるプロセス300のフローチャートである。リンクOTFプロセス300は、情報のフィードを定常的に評価し、スコア付けし、リンクし、そしてクラスター化することができる。リンクOTFサブモジュールは、複数のアルゴリズムを使用してレコードリンケージ206を遂行する。ステップ204の候補サーチ結果は、リンクOTFモジュール300へ定常的にフィードされる。データの入力に続いて、一致スコア付けアルゴリズムが適用され（ステップ302）、ここでは、1つ以上の一致スコア付けアルゴリズムがMemDBの複数のサーチノードに同時に適用される一方、とりわけ、ストリンク編集距離、表音及び意味のような複数の特徴属性を考慮して、関連する結果を評価及びスコア付けするために曖昧キーサーチを遂行する。

#### 【0056】

その後、一致スコア付けアルゴリズム適用ステップ302の間に識別された全ての候補レコードを互いに比較するために、リンキングアルゴリズムの適用ステップ304が追加

10

20

30

40

50

される。リンキングアルゴリズムの適用 3 0 4 は、M e m D B の複数のサーチノードの内部で遂行される曖昧キーサーチのスコア付けされた結果をフィルタリング及び評価できる 1 つ以上の分析リンキングアルゴリズムの使用を含む。ある例では、M e m D B における識別された候補レコードの集合体にわたる 2 つ以上の特徴の共起を分析して、プロセスの精度を改善する。リンキングアルゴリズムの適用 3 0 4 には、異なる特徴属性に関連した異なる重み付けモデル及び信頼性スコアが考慮される。

#### 【 0 0 5 7 】

リンキングアルゴリズムの適用ステップ 3 0 4 の後に、リンクされた結果が関連特徴のクラスターに配置され、そしてステップ 3 0 6 において、リンクされたレコードのクラスターの返送の一部分として返送される。

10

#### 【 0 0 5 8 】

図 4 は、図 1 を参照して上述した非構造化テキストにおいて特徴を曖昧性除去するシステム 4 0 0 の一実施形態を例示する図である。このシステム 4 0 0 は、インメモリデータベースをホストし、そして 1 つ以上のノードを含む。

#### 【 0 0 5 9 】

一実施形態によれば、システム 4 0 0 は、1 つ以上のドキュメント内の特徴を曖昧性除去するため複数の特殊目的コンピュータモジュール 4 0 1、4 0 2、4 1 1、4 1 2 及び 4 1 4 (以下に述べる) のコンピュータインストラクションを実行する 1 つ以上のプロセッサを備えている。図 4 に示すように、ドキュメント入力モジュール 4 0 1、4 0 2 は、インターネットベースのソース及び/又はドキュメントの生のコーパスからドキュメントを受け取る。多数の新たなドキュメントがネットワーク接続 4 0 4 を通してドキュメント入力モジュール 4 0 2 へアップロードされる。それ故、ソースは、常時、新たな知識を得て、ユーザワークステーション 4 0 6 により更新され、そのような新たな知識は、ステイックな仕方ですべてリンクされない。従って、評価されるべきドキュメントの数は、無限に増加する。

20

#### 【 0 0 6 0 】

この評価は、M e m D B コンピュータ 4 0 8 を経て達成される。M e m D B 4 0 8 は、高速の曖昧性除去プロセスを促進し、曖昧性除去プロセスをオンザフライで促進し、これは、M e m D B 4 0 8 に貢献しようとする最新情報の受信を促進する。特徴をリンクするための種々の方法が使用され、これは、重み付けされたモデルを本質的に使用して、どのエンティティタイプが最も重要であるか決定し、どれがより大きな重みを有するか決定し、そして信頼性スコアに基づき、正しい特徴の抽出及び曖昧性除去がどれほどの信頼性で遂行されたか決定し、且つ正しい特徴が結果の特徴クラスターに向かうことを決定する。図 4 に示すように、より多くのシステムノードが並列に機能するほど、プロセスは、より効率的となる。

30

#### 【 0 0 6 1 】

種々の実施形態によれば、新たなドキュメントがドキュメント入力モジュール 4 0 1、4 0 2 を経てネットワーク接続 4 0 4 を通してシステム 4 0 0 に到着するとき、特徴抽出が抽出モジュール 4 1 1 を経て遂行され、次いで、特徴の曖昧性除去が新たなドキュメントにおいて M e m D B 4 0 8 の特徴曖昧性除去サブモジュール 4 1 4 を経て遂行される。ある実施形態では、新たなドキュメントの特徴曖昧性除去が遂行された後に、抽出された新たな特徴 4 1 0 は、リンク O T F サブモジュール 4 1 2 を通過するために M e m D B に含まれ、ここで、特徴は、比較され及びリンクされ、そして曖昧性除去された特徴 1 1 0 の特徴 I D が質問からの結果としてユーザに返送される。特徴 I D に加えて、曖昧性除去された特徴を定義する結果の特徴クラスターが任意に返送されてもよい。

40

#### 【 0 0 6 2 】

M e m D B コンピュータ 4 0 8 は、装置メインメモリにデータレコードを記憶するように構成されたデータベース管理システム ( D B M S ) (図示せず) により制御されるレコードにデータを記憶するデータベースであり、これは、データを「ディスク」メモリに記憶する従来のデータベース及び D B M S モジュールと対照的である。従来のデ

50

ィスクストレージは、装置のハードディスクへの読み取り及び書き込みコマンドをプロセッサ（ＣＰＵ）が実行することを要求し、従って、ＣＰＵがデータのためのメモリ位置を位置付け（即ち、シークし）及び検索するインストラクションを実行した後に、そのメモリ位置におけるデータとのある形式のオペレーションを遂行することを要求する。インメモリデータベースシステムは、メインメモリに入れられて適宜にアドレスされるデータにアクセスし、従って、ＣＰＵにより遂行されるインストラクションの数を軽減し、そしてハードディスクのデータをＣＰＵがシークするのに関連したシークタイムを排除する。

【 0 0 6 3 】

インメモリデータベースは、ノードの各リソース（例えば、メモリ、ディスク、プロセッサ）をアグリゲートするように構成された１つ以上のノードを含むコンピューティングシステムである分散型コンピューティングアーキテクチャーにおいて実施される。ここに開示されるように、インメモリデータベースをホストするコンピューティングシステムの実施形態は、１つ以上のノードの間でデータベースのデータレコードを分散しそして記憶する。ある実施形態では、これらのノードは、ノードの「クラスター」へと形成される。ある実施形態では、ノードのこれらクラスターは、データベース情報の部分又は「集合体」を記憶する。

【 0 0 6 4 】

種々の実施形態は、共起トピック、キーフレーズ、接近用語、イベント、ファクト及びトレンド人気指数のような二次特徴を記憶するように構成された進化する効率的にリンク可能な特徴知識ベースを使用するコンピュータ実行の特徴曖昧性除去技術を提供する。ここに開示する実施形態は、知識ベースに記憶された特徴に対して所与の抽出特徴を分析する上で役立つ関連二次特徴の次元に基づいて簡単な概念的距離尺度から精巧なグラフクラスター化解決策まで変化し得る種々様々なリンキングアルゴリズムを経て遂行される。加えて、それらの実施形態は、既存の特徴エントリの二次特徴を更新するだけでなく、知識ベースに追加できる新たな特徴を発見することでそれを拡張もする能力により既存の特徴知識ベースを進化させる解決策を導入することができる。

【 0 0 6 5 】

曖昧性除去解決策の実施形態は、トピックモデリング解決策を使用して、トピック推論としてモデリングされる自動重み付け（全ての二次特徴にわたる）リンキングプロセスを提供する。この自動重み付け型リンキングプロセスをサポートするため、それら実施形態は、多数のコンポーネント（二次特徴）を条件独立としてサポートできるマルチコンポーネントＬＤＡ（ＭＣ－ＬＤＡ）と称される新規なトピックモデリング解決策を構築するように従来のＬＤＡトピックモデリングを拡張する。又、モデリング解決策の実施形態は、トレーニング中にコンポーネントの重みを自動的に学習し、そしてそれを曖昧性除去に関する推論（リンキング）のために使用することができる。曖昧性除去のために導入されるＭＣ－ＬＤＡ解決策は、曖昧性除去精度を高めるために導入できる付加的な数の二次特徴のためにスケーリングすることができる。

【 0 0 6 6 】

図５は、上述した図４のシステム４００によって使用されるマルチコンポーネント条件独立のレイテントディリクレアロケーション（ＭＣ－ＬＤＡ）トピックコンピュータモデリング解決策の実施形態のグラフィック表示である。ここに示す実施形態では、各コンポーネントブロックは、例えば、図５に示すパラメータで初期化される図４のＭｅｍＤＢ４０８を経て実行される、知識ベースにわたる各二次特徴のモデリングを表す。

【 0 0 6 7 】

図６は、上述した図５に使用されるＭＣ－ＬＤＡトピックモデルのギブスサンプリング方程式の実施形態を示す。このサンプリング解決策の実施形態は、個々のコンポーネント（二次特徴）の重みを自動的に且つ効率的にトレーニングする上で図４のシステム４００の助けとなる。

【 0 0 6 8 】

図７は、例えば、図７に示すパラメータで初期化される図４のシステム４００のＭｅｍ

10

20

30

40

50

DB 408を経て実行される、図5-6のMC-LDAトピックモデルにおけるトレーニング及び推論のための確率論的变化推論アルゴリズムのコンピュータ実行の実施形態を示す。この推論方法の実施形態は、全ての二次特徴（当該ドキュメントから抽出された）を入力として取り上げそして重み付けされたトピックを出力として与えることにより、リンク性/曖昧性除去プロセスをトピック推論としてモデリングするように容易に適用される。これらの重み付けされたトピックは、次いで、記憶された特徴知識ベースエントリに対して類似性スコアを計算するのに使用できる。

#### 【0069】

図8は、MC-LDAトピックモデルに対するサンプルトピックを示すテーブルである。図8は、一実施形態により、例えば、図4のシステム400のMemDB 408を経て実行される、モデルの各コンポーネントに対するトップスコア付け表面フォームを示す。

#### 【0070】

例#1は、当該特徴（一次特徴）がフットボール選手のJohn DoeでありそしてユーザがJohn Doeについて言及するニュースの監視を希望する場合に、非構造化テキストにおける特徴を曖昧性除去する方法100を適用するものである。ある実施形態によれば、John Doeについて述べるドキュメント入力102がネットワークにアップロードされる。ドキュメント入力102の特徴が抽出されて、MemDB 408に含まれ、曖昧性除去されて、一次特徴（John Doe）に関連した二次特徴のクラスターにリンクされ、そして同様の特徴の既存のクラスターと比較される。方法100は、異なる特徴ID及び特徴IDの関連クラスターを出力し、これは、John Doeに対する全ての関連二次特徴、例えば、エンジニアのJohn Doe；教師のJohn Doe；及びフットボール選手のJohn Doe；を含む。同様の二次特徴を伴う他の一次特徴、例えば、ニックネーム又は省略名が考えられる。フットボール選手のJohn Doeと同じチームから、同じ年齢及び経験の「JD」フットボール選手は、同じ一次特徴と考えられる。それ故、フットボール選手のJohn Doeに関連した全てのドキュメントは、容易にアクセスすることができる。

#### 【0071】

例#2は、一次特徴が画像である場合に、非構造化ドキュメントにおける特徴を曖昧性除去する方法100を適用するものである。ある実施形態によれば、方法100は、特徴の抽出104を含み、ここで、特徴は、とりわけ、縁及び形状のような一般的な属性であるか、或いはとりわけ、タンク、個人及び時計のような特定の属性である。例えば、新たな画像が入力され、ここで、画像は、特定の形状（例えば、方形、個人又は車の形状）のような二次特徴を有し、二次特徴が抽出されてMemDB 408に含まれ、ここで、同様の二次特徴を有する他の全ての画像の間で一致が見出される。ここに示す実施形態によれば、特徴は、画像のみを含み、即ちテキストは、特徴として含まれない。

#### 【0072】

例#3は、一次特徴がイベントである場合に、非構造化テキストにおける特徴を曖昧性除去する方法100を適用するものである。ある実施形態によれば、質問がなされたときに、方法100は、ユーザが、とりわけ、地震、火災、又は伝染病の発生のようなイベントに関連した結果を受け取ることができるようにする。方法100は、特徴の抽出104及び特徴の曖昧性除去108を遂行して、イベントに関連した特徴を見出すと共に、曖昧性除去された特徴110の特徴IDを与える。

#### 【0073】

例#4は、1つ以上のイベントの発生の予想がなされる場合の方法100の実施形態である。ある実施形態によれば、ユーザは、オペレーションの前に当該特徴及びイベントを前もって指示し、それ故、当該イベントに関連した異なる特徴間のリンクが前もって確立される。関連特徴が高い発生数でネットワークに現われるとき、方法100は、関連特徴の発生数増加に基づいて、当該事象が発生することを予想する。切迫したイベントが検出されると、ユーザに警報が送られる。例えば、タイからの保健省に対して仕事をするユーザは、デング熱の伝染病発生についての警報を受け取ることを選択する。例えば、ソーシ

10

20

30

40

50

ャルネットワークからの他のユーザ406がデング熱の兆候又は包括を含めたコメントを病院へアップロードするとき、方法100は、ソーシャルネットワークからの全ての関連コメントを曖昧性除去し、そして関連情報を含めたユーザ406の数を考慮して、デング熱の伝染病発生が生じることを予想し、保健省の職員に警報する。それ故、保健省の職員は、付加的な形跡を得て、影響のある共同体への更なる対策を取り、伝染病が広がらないようにする。

【0074】

例#5は、一次特徴が地理的な場所の名前である場合の方法100の適用である。一実施形態によれば、方法100は、都市の名前を曖昧性除去するのに使用され、曖昧性除去サブモジュールにおいて二次特徴に異なるスコア付け重みが関連付けられる。例えば、方法100は、Paris、TexasをParis、Franceから曖昧性除去するのに使用される。

10

【0075】

例#6は、一次特徴が、とりわけ、個人、イベント、又は会社に関連した感情であり、その感情が、とりわけ、個人、イベント、又は会社に関する肯定的又は否定的コメントであって、ソーシャルネットワークを含む適当なソースから供給される場合に、非構造化テキストにおける特徴を曖昧性除去する方法100を適用するものである。ある実施形態によれば、方法100は、会社が一般大衆の中で有している容認性を確認するために使用される。

【0076】

20

例#7は、特徴の信頼性スコアを高めるために人間の確認を含む方法100の実施形態である。ある実施形態によれば、リンクOTFプロセス300(図4)は、ユーザにより支援され、ユーザは、曖昧性除去された特徴が正しく曖昧性除去されたかどうか指示し、そして2つの異なるクラスターが1つでなければならないかどうか指示し、これは、ユーザが知っている2つの異なる一次特徴が同じであるときに方法100が(全ての特徴及びトピック共起情報を考慮して)何を指示するかを意味する。それ故、そのクラスターに関連した信頼性スコアが高くなり、従って、特徴が正しく曖昧性除去されたという確率が高くなる。

【0077】

例#8は、曖昧性除去プロセス200及びリンクOTFプロセス300を使用する方法100の実施形態である。この例では、リンキングアルゴリズムの適用304に使用されるリンキングアルゴリズムは、1000msの期間内に0.85より高い信頼性スコアを与えるように構成される。

30

【0078】

例#9は、曖昧性除去プロセス200及びリンクOTFプロセス300を使用する方法100の実施形態である。この例では、リンキングアルゴリズムの適用304に使用されるリンキングアルゴリズムは、300ms以下の期間内に0.80より高い信頼性スコアを与えるように構成される。この例に使用されるアルゴリズムは、例#8に使用されるアルゴリズムに比して短い期間内に応答を与えるが、一般的に、低い信頼性スコアを返送する。

40

【0079】

例#10は、曖昧性除去プロセス200及びリンクOTFプロセス300を使用する方法100の実施形態である。この例では、リンキングアルゴリズムの適用304に使用されるリンキングアルゴリズムは、一般的に3000msを越える期間内に0.90より高い信頼性スコアを与えるように構成される。この例に使用されるアルゴリズムは、例#8に使用されるアルゴリズムにより返送されるものより一般的に大きな信頼性スコアをもつ応答を与えるが、著しく長い期間を一般的に要求する。

【0080】

例#11は、複数のソースからのドキュメントの大きなコーパスにおいてeディスカバリーを遂行するために非構造化テキストにおける特徴を曖昧性除去する方法100の一例

50

である。複数のリソースからのドキュメントの大きなコーパスが与えられると、それらドキュメントにおける全ての特徴を曖昧性除去するための方法100の適用は、コーパスにおいて全ての特徴を発見できるようにする。発見された特徴の集合体は、特徴に関連した全てのドキュメントの発見及び関連特徴の発見に更に使用することができる。

【0081】

以上の方法の説明及びプロセスフロー図は、単なる例示として示されたもので、種々の実施形態のステップを、提示した順序で遂行しなければならないことを要求し又は意味することは意図されない。当業者に明らかなように、前記実施形態におけるステップは、任意の順序で遂行されてもよい。「次いで(then)」、「次に(next)」、等のワードは、ステップの順序を限定するものではなく、これらのワードは、単に、方法の説明を通して読者を誘導するのに使用されるだけである。プロセスフロー図は、オペレーションを一連のプロセスとして示すが、多数のオペレーションを並列に又は同時に遂行することもできる。加えて、オペレーションの順序は、再構成してもよい。プロセスは、方法、機能、手順、サブルーチン、サブプログラム、等に対応する。プロセスが機能に対応するとき、その終了は、コーリング機能又はメイン機能への機能の復帰に対応する。

【0082】

ここに開示する実施形態に関連して述べた種々の例示的論理ブロック、モジュール、回路及びアルゴリズムステップは、電子的ハードウェア、コンピュータソフトウェア又はその両方の組み合わせとして具現化されてもよい。ハードウェア及びソフトウェアのこの互換性を明確に示すために、種々の例示的コンポーネント、ブロック、モジュール、回路、及びステップは、それらの機能に関して一般的に説明された。そのような機能がハードウェアとして具現化されるかソフトウェアとして具現化されるかは、システム全体に課せられる特定アプリケーション及び設計上の制約に依存する。当業者であれば、ここに述べた機能を特定アプリケーションごとに色々な仕方で具現化できるが、そのような具現化の判断は、本発明の範囲から逸脱すると解釈されてはならない。

【0083】

コンピュータソフトウェアで具現化される実施形態は、ソフトウェア、ファームウェア、ミドルウェア、マイクロコード、ハードウェア記述言語、又はその組み合わせで具現化される。コードセグメント又はマシン実行可能なインストラクションは、手順、機能、サブプログラム、プログラム、ルーチン、サブルーチン、モジュール、ソフトウェアパッケージ、クラス、或いはインストラクション、データ構造体又はプログラムステートメントの組合せを表わす。コードセグメントは、情報、データ、アークギュメント、パラメータ又はメモリコンテンツを通し及び/又は受け取ることにより別のコードセグメント又はハードウェア回路に結合される。情報、アークギュメント、パラメータ、データ、等は、メモリ共有、メッセージ通過、トークン通過、ネットワーク送信、等を含む適当な手段を経て通され、転送され又は送信される。

【0084】

これらのシステム及び方法を実施するのに使用される実際のソフトウェアコード又は特殊な制御ハードウェアは、本発明を限定するものではない。従って、システム及び方法のオペレーション及び振舞いは、ここでの記載に基づいてシステム及び方法を実施するようにソフトウェア及び制御ハードウェアを設計できることを理解して、特定のソフトウェアコードを参照せずに説明した。

【0085】

ソフトウェアで実施されるときに、機能は、非一時的コンピュータ読み取り可能な又はプロセッサ読み取り可能なストレージ媒体に1つ以上のインストラクション又はコードとして記憶される。ここに開示する方法又はアルゴリズムのステップは、コンピュータ読み取り可能な又はプロセッサ読み取り可能なストレージ媒体に存在するプロセッサ実行可能なソフトウェアモジュールにおいて実施される。非一時的なコンピュータ読み取り可能な又はプロセッサ読み取り可能な媒体は、ある場所から別の場所へのコンピュータプログラムの転送を容易にするコンピュータストレージ媒体及び有形のストレージ媒体の両方を含

10

20

30

40

50

む。非一時的なプロセッサ読み取り可能なストレージ媒体は、コンピュータによりアクセスされる利用可能な媒体である。これに限定されないが、一例として、そのような非一時的なプロセッサ読み取り可能な媒体は、RAM、ROM、EEPROM、CD-ROM又は他の光学ディスクストレージ、磁気ディスクストレージ又は他の磁気ストレージ装置、或いはインストラクション又はデータ構造体の形態で望ましいプログラムコードを記憶するのに使用され且つコンピュータ又はプロセッサによりアクセスされる他の有形のストレージ媒体を含む。ここで使用するディスク(disk & disc)とは、コンパクトディスク(CD)、レーザーディスク(登録商標)、光学ディスク、デジタル多様性ディスク(DVD)、フロッピーディスク、及びブルーレイディスクを含み、ここで、ディスク(disk)は、通常、データを磁氣的に再生するものであり、一方、ディスク(disc)は、データをレーザーで光学的に再生するものである。前記の組み合わせも、コンピュータ読み取り可能な媒体の範囲内に包含される。加えて、方法又はアルゴリズムのオペレーションは、コンピュータプログラム製品に合体される非一時的プロセッサ読み取り可能な媒体及び/又はコンピュータ読み取り可能な媒体にコード及び/又はインストラクションの1つ又は組み合わせ或いはセットとして存在する。

10

#### 【0086】

技術の種々のコンポーネントは、分散型ネットワーク及び/又はインターネットの遠隔部分に、或いは専用のセキュア、アンセキュア及び/又は暗号化システム内に配置できることが明らかである。従って、システムのコンポーネントは、1つ以上の装置に結合するか、又はテレコミュニケーションネットワークのような分散型ネットワークの特定ノードに共通配置できることが明らかである。以上の説明から明らかなように、計算効率の理由で、システムのコンポーネントは、システムのオペレーションに影響することなく、分散型ネットワーク内の任意の位置に配置することができる。更に、それらのコンポーネントは、専用マシンに埋め込むこともできる。

20

#### 【0087】

更に、エレメントを接続する種々のリンクは、ワイヤード又はワイヤレスリンク又はその組み合わせ、或いは接続されたエレメントへ及びそこからデータを供給及び/又は通信することのできる他の既知の又は今後開発されるエレメントであることが明らかである。ここで使用するモジュールという語は、エレメントに関連した機能を遂行できる既知の又は今後開発されるハードウェア、ソフトウェア、ファームウェア、又はその組み合わせを指す。又、ここで使用する決定、計算及びコンピューティング、並びにその変形の語は、交換可能に使用され、そして任意のタイプの方法、プロセス、数学演算又は技術を包含する。

30

#### 【0088】

ここに開示する実施形態の前記説明は、当業者が本発明を実施又は利用できるようにするためになされたものである。これら実施形態に対する種々の変更は、当業者に容易に明らかであり、そしてここに定義する一般的な原理は、本発明の精神又は範囲から逸脱せずに他の実施形態に適用される。従って、本発明は、ここに示す実施形態に限定されるものではなく、特許請求の範囲並びにここに開示した原理及び新規な特徴に一致する最も広い範囲と調和されるべきである。

40

#### 【0089】

以上に述べた実施形態は、例示に過ぎない。当業者であれば、ここに述べた特定例に対して置き換えられ且つ依然として本発明の範囲内に入る多数の代替的コンポーネント及び実施形態が認識されよう。

#### 【符号の説明】

#### 【0090】

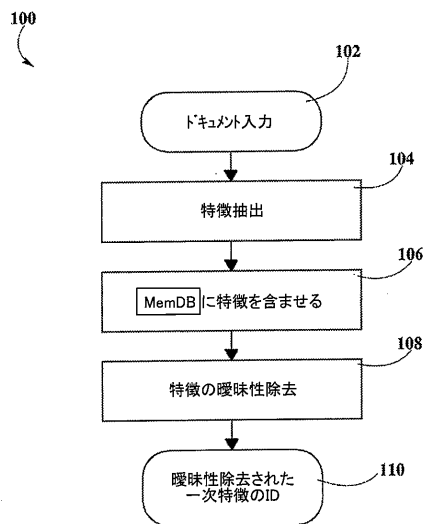
- 400：システム
- 401、402：ドキュメント入力モジュール
- 404：ネットワーク接続
- 406：ユーザワークステーション

50

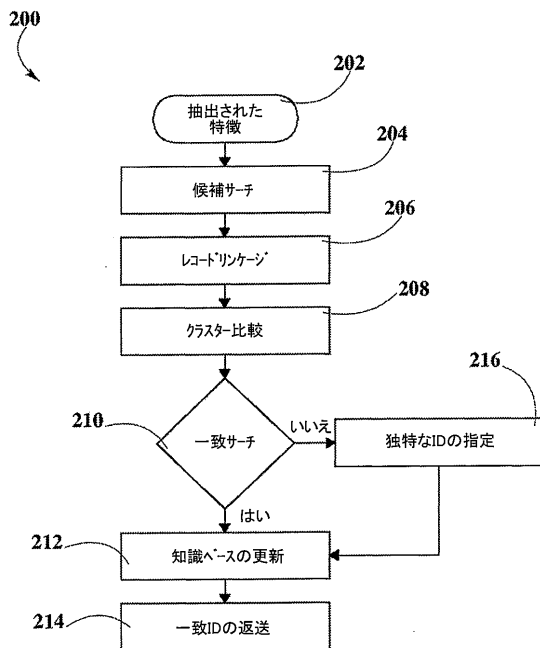


408 : MemDB コンピュータ  
410 : 抽出された新たな特徴  
411 : 抽出ノード  
412 : リンクOTFサブモジュール

【図1】



【図2】



【図 3】

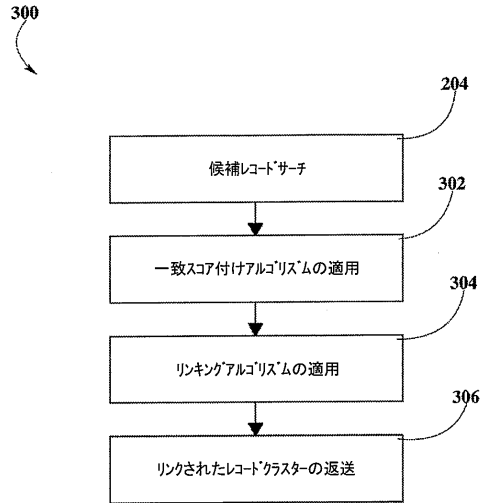


FIG. 3

【図 4】

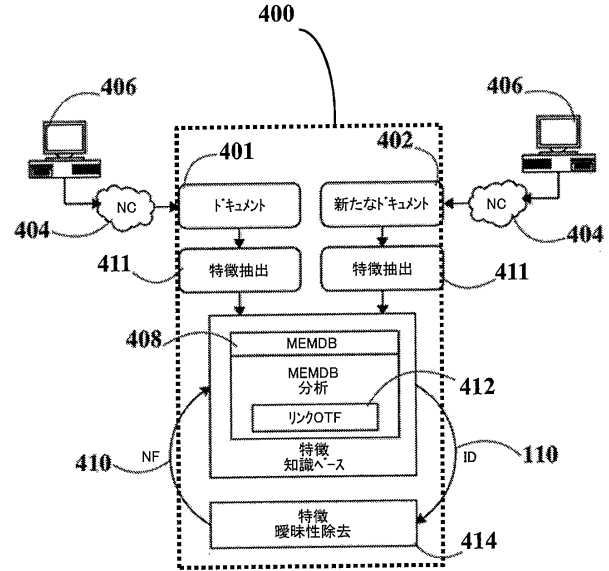
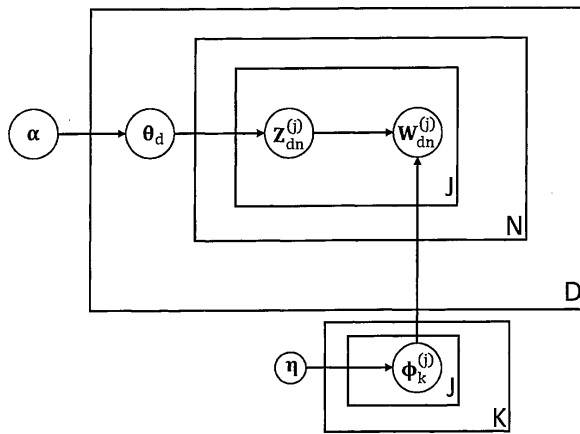


FIG. 4

【図 5】



- $J$  = マルチドキュメントコンポーネントの数  
 $V^{(j)}$  =  $j$  番目のコンポーネントの語彙における用語の数  
 $D$  = ドキュメントの数  
 $N$  = ドキュメントの(コンポーネント)におけるトークンの数(実際には  $j$  及び  $d$  の両方に依存する)  
 $K$  = トピックの数  
 $\alpha$  = 混合比率でのハイパーパラメータ(対称的である場合には  $K$  ベクトル又はスカラー)  
 $\eta^{(j)}$  = 混合コンポーネントでのハイパーパラメータ(対称的である場合には  $V^{(j)}$  ベクトル又はスカラー)  
 $\theta_d$  = ドキュメント  $d$  に対するトピック混合比率  
 $\phi_k^{(j)}$  =  $k$  番目のトピックの  $j$  番目のコンポーネントに対する混合コンポーネント  
 $z_{dn}^{(j)}$  = ドキュメント  $d$  の  $j$  番目のコンポーネントにおける  $n$  番目のワードに対するトピックを選択する混合インジケータ  
 $w_{dn}^{(j)}$  = ドキュメント  $d$  の  $j$  番目のコンポーネントにおける  $n$  番目のワードに対する項目インジケータ

FIG. 5

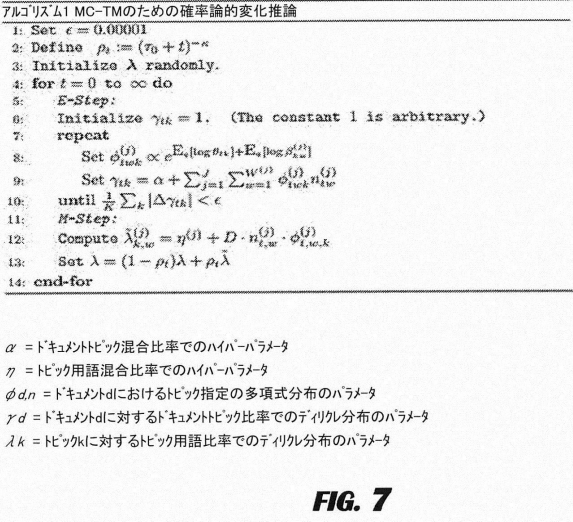
【図 6】

$$\begin{aligned}
 p(z_i^{(j)} = k | w_i^{(j)} = t_i, z_{-i}^{(j)}, w_{-i}^{(j)}, w^{(-j)}, \alpha, \eta) &\propto \frac{n_{k,i-1}^{(j,k)} + \alpha}{\sum_{k'=1}^K (n_{k',i-1}^{(j,k')} + \alpha)} \cdot \frac{n_{k,i-1}^{(j,i)} + \beta^{(j)}}{\sum_{i'=1}^{V^{(j)}} (n_{k,i-1}^{(j,i')} + \beta^{(j)})} \\
 E\{\phi_{k,i}^{(j)} | z_{-i}^{(j)}, w_{-i}^{(j)}, \beta\} &= \frac{n_k^{(j,i)} + \beta^{(j)}}{\sum_{i'=1}^{V^{(j)}} (n_k^{(j,i')} + \beta^{(j)})} \\
 E\{\theta_{d,k} | z^{(1)}, \dots, z^{(j)}, \alpha\} &= \frac{\sum_{j=1}^J (n_{d,k}^{(j,k)} + \alpha)}{\sum_{k=1}^K \sum_{j=1}^J (n_{d,k}^{(j,k)} + \alpha)}
 \end{aligned}$$

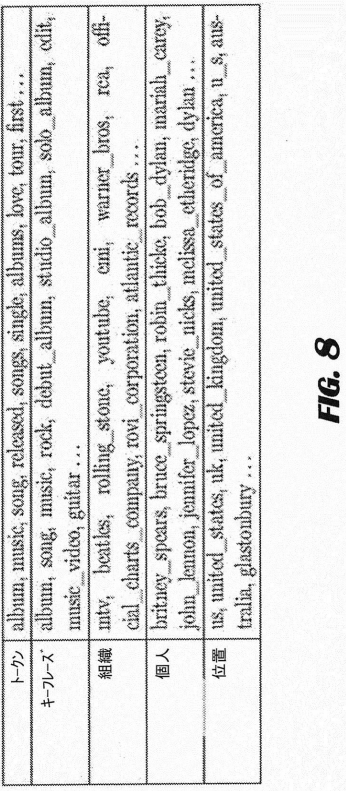
- $J$  = マルチドキュメントコンポーネントの数  
 $V^{(j)}$  =  $j$  番目のコンポーネントの語彙における用語の数  
 $K$  = トピックの数  
 $\alpha$  = 混合比率でのハイパーパラメータ(対称的である場合には  $K$  ベクトル又はスカラー)  
 $\eta$  = 混合コンポーネントでのハイパーパラメータ(実際には、各々、長さ  $V^{(j)}$  の  $J$  個のベクトル)  
 $\theta_d$  = ドキュメント  $d$  に対するトピック混合比率  
 $\phi_k^{(j)}$  =  $k$  番目のトピックの  $j$  番目のコンポーネントに対する混合コンポーネント  
 $z_{dn}^{(j)}$  = ドキュメント  $d$  の  $j$  番目のコンポーネントにおける  $n$  番目のワードに対するトピックを選択する混合インジケータ  
 $w_{dn}^{(j)}$  = ドキュメント  $d$  の  $j$  番目のコンポーネントにおける  $n$  番目のワードに対する項目インジケータ  
 $n_{k,i}^{(j,i)}$  = トピック  $k$  に指定されたドキュメント  $d$  のコンポーネント  $j$  における用語の数  
 $n_k^{(j,i)}$  = コーパスのコンポーネント  $j$  におけるトピック  $k$  に用語  $i$  が指定された回数

FIG. 6

【図 7】



【図 8】



## フロントページの続き

- (74)代理人 100109070  
弁理士 須田 洋之
- (74)代理人 100109335  
弁理士 上杉 浩
- (74)代理人 100120525  
弁理士 近藤 直樹
- (74)代理人 100139712  
弁理士 那須 威夫
- (72)発明者 ライトナー スコット  
アメリカ合衆国 ヴァージニア州 20175 リーズバーグ レッドヒル マナー コート 2  
2596
- (72)発明者 ウェックザー フランツ  
アメリカ合衆国 オハイオ州 45370 スプリング ヴァレー イースト センターヴィル  
ロード 3942
- (72)発明者 ボッツ サンジェイ  
アメリカ合衆国 オハイオ州 45440 デイトン サマーセット パス 4408
- (72)発明者 デイヴ ラケシュ  
アメリカ合衆国 オハイオ州 45440 デイトン アマースト ベンド 376
- (72)発明者 フラッグ ロバート  
アメリカ合衆国 メイン州 04101 ポートランド イースタン ブロムナード 6 アパー  
トメント 1

審査官 齊藤 貴孝

- (56)参考文献 米国特許出願公開第2002/0165847(US,A1)  
米国特許出願公開第2013/0290665(US,A1)  
特開2013-239162(JP,A)  
特開2003-150442(JP,A)

- (58)調査した分野(Int.Cl., DB名)  
G06F 17/30