

(19) 日本国特許庁(JP)

(12) 公開特許公報(A)

(11) 特許出願公開番号

特開2008-9671

(P2008-9671A)

(43) 公開日 平成20年1月17日(2008.1.17)

| | | |
|-----------------------------|-----------------|-------------|
| (51) Int. Cl. | F I | テーマコード (参考) |
| G06F 17/30 (2006.01) | G06F 17/30 210D | 5B075 |
| | G06F 17/30 380E | |
| | G06F 17/30 170A | |

審査請求 未請求 請求項の数 16 O L (全 24 頁)

| | | | |
|-----------|------------------------------|----------|---|
| (21) 出願番号 | 特願2006-178922 (P2006-178922) | (71) 出願人 | 301022471 独立行政法人情報通信研究機構 東京都小金井市貫井北町4-2-1 |
| (22) 出願日 | 平成18年6月29日 (2006.6.29) | (74) 代理人 | 100119161 弁理士 重久 啓子 |
| | | (74) 代理人 | 100121511 弁理士 小田 直 |
| | | (72) 発明者 | 村田 真樹 東京都小金井市貫井北町4-2-1 独立 行政法人情報通信研究機構内 |
| | | Fターム(参考) | 5B075 ND03 NK02 NK32 NR12 NR15 NS03 PQ02 PQ15 PQ46 PR04 QP01 UU06 |

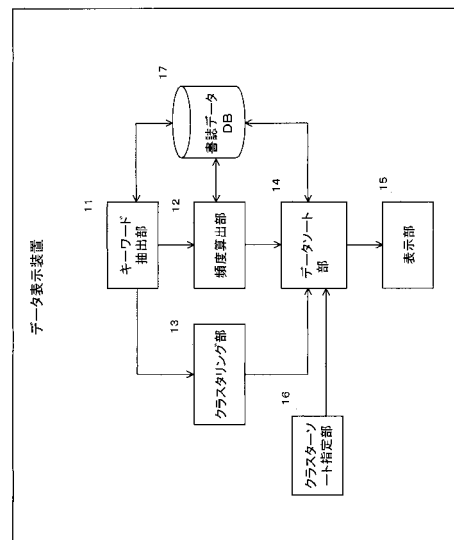
(54) 【発明の名称】 データ表示装置、データ表示方法及びデータ表示プログラム

(57) 【要約】

【課題】 文書の大凡の内容を一見して把握できるような表示を行うことを可能とする。

【解決手段】 キーワード抽出部11が書誌データDB17に蓄積されている文書群に含まれるキーワードを抽出し、頻度算出部12がキーワード抽出部11によって抽出された各キーワードの、書誌データDB17に蓄積されている文書群中に出現した頻度を算出し、データソート部14が書誌データDB17に蓄積されている各文書の文書データを抽出し、各キーワードの頻度がより高いキーワードを多く含む文書の順に、各文書の文書データをソートし、クラスタリング部13がキーワード抽出部11によって抽出された各キーワードをクラスターにクラスタリングし、表示部15がデータソート部14によってソートされた文書データに対応付けて、上記クラスター中のキーワードのうち当該ソートされた文書データが持つ文書が含むキーワードを、クラスター毎に画面表示する。

【選択図】 図1



【特許請求の範囲】**【請求項 1】**

データ表示装置であって、

複数の文書から構成される文書群中に含まれるキーワードを抽出するキーワード抽出手段と、

前記抽出されたキーワードを各キーワードが属するクラスターにクラスタリングするクラスタリング手段と、

前記クラスターに属するキーワードのうちの前記文書群を構成する文書に含まれるキーワード、又は、各文書中における前記クラスターに属するキーワードの有無を示す情報を、前記各文書と前記クラスターとに対応付けて画面表示する表示手段とを備える

10

ことを特徴とするデータ表示装置。

【請求項 2】

請求項 1 に記載のデータ表示装置において、

前記クラスタリング手段は、前記各キーワードのベクトル空間上の位置を示す位置ベクトルを生成し、生成された前記位置ベクトルが示す前記各キーワードの位置情報に基づいて、前記各キーワードが属するクラスターを決定する

ことを特徴とするデータ表示装置。

【請求項 3】

請求項 1 又は請求項 2 に記載のデータ表示装置において、

前記表示手段は、予め定められた順序に基づいて、前記各クラスターを並び替えて画面表示する

20

ことを特徴とするデータ表示装置。

【請求項 4】

請求項 1 又は請求項 2 に記載のデータ表示装置において、

前記表示手段は、各クラスターに属するキーワード又は前記各クラスターに属するキーワードが出現する文書の分布に基づいて、前記各クラスターを並び替えて画面表示する

ことを特徴とするデータ表示装置。

【請求項 5】

請求項 4 に記載のデータ表示装置において、

前記表示手段は、各クラスターに属するキーワードが出現する文書数について昇順又は降順に前記各クラスターを並び替えて画面表示する

30

ことを特徴とするデータ表示装置。

【請求項 6】

請求項 2 に記載のデータ表示装置において、

前記クラスタリング手段は、前記各キーワードの位置ベクトルの平均を前記各キーワードが属するクラスターの位置ベクトルとし、各クラスターの位置ベクトルに基づいて、クラスター同士の距離を求め、求めた距離が近いクラスター同士が近い位置に並び、求めた距離が遠いクラスター同士が離れた位置に並ぶように、前記各クラスターを並び替え、

前記表示手段は、前記並び替えられた各クラスターに属するキーワードのうちの前記文書群を構成する文書に含まれるキーワード、又は、各文書中における前記並び替えられた各クラスターに属するキーワードの有無を示す情報を、前記各文書と前記並び替えられた各クラスターとに対応付けて画面表示する

40

ことを特徴とするデータ表示装置。

【請求項 7】

請求項 6 に記載のデータ表示装置において、

前記クラスタリング手段は、ベクトル空間上に並んで配置された複数の基準の位置ベクトルを予め定義し、前記各クラスターの位置ベクトルを、前記各クラスターの位置ベクトルとの距離が最も近い基準の位置ベクトルの近くに配置することにより、前記各クラスターを並び替える

ことを特徴とするデータ表示装置。

50

【請求項 8】

請求項 6 又は請求項 7 に記載のデータ表示装置において、

前記クラスタリング手段は、前記各文書がどのクラスターのどの単語を何個含んでいるかの情報を求め、求めた情報に基づいて、前記各文書のベクトル空間上の位置を示す位置ベクトルを求め、求めた前記各文書の位置ベクトルに基づいて、文書同士の距離を求め、求めた距離が近い文書同士が近い位置に並び、求めた距離が遠い文書同士が離れた位置に並ぶように、各文書を並び替える

ことを特徴とするデータ表示装置。

【請求項 9】

請求項 4 乃至請求項 8 のいずれか 1 項に記載のデータ表示装置において、

前記表示手段は、並び順序がより上位のクラスターに属するキーワードを含む文書の順に、前記文書に含まれるキーワード、又は、前記文書中における前記クラスターに属するキーワードの有無を示す情報を画面表示する

ことを特徴とするデータ表示装置。

【請求項 10】

請求項 1 乃至請求項 9 のいずれか 1 項に記載のデータ表示装置において、

前記表示手段は、前記クラスタリング手段によってクラスタリングされたクラスターを選択し、選択されたクラスターに対応付けて、前記クラスターに属するキーワードのうちの前記文書群を構成する文書に含まれるキーワード、又は、各文書中における前記クラスターに属するキーワードの有無を示す情報を画面表示する

ことを特徴とするデータ表示装置。

【請求項 11】

請求項 1 に記載のデータ表示装置において、

前記クラスタリング手段は、前記各キーワードのベクトル空間上の位置を示す位置ベクトルを生成し、生成された前記位置ベクトルが示す前記各キーワードの位置情報に基づいて、キーワード同士の距離を求め、求めた距離が近いキーワード同士が近い位置に並び、求めた距離が遠いキーワード同士が離れた位置に並ぶように、各キーワードを並び替え、

前記表示手段は、並び替えられたキーワードのうちの前記文書群を構成する文書に含まれるキーワード、又は、各文書中における前記並び替えられたキーワードの有無を示す情報を、前記各文書に対応付けて画面表示する

ことを特徴とするデータ表示装置。

【請求項 12】

請求項 1 に記載のデータ表示装置において、

前記表示手段は、各クラスターを画面上に並ばせる順序を選択し、前記選択された順序がより上位のクラスターに属するキーワードを含む文書の順に、前記文書に含まれるキーワード、又は、前記文書中における前記クラスターに属するキーワードの有無を示す情報を画面表示する

ことを特徴とするデータ表示装置。

【請求項 13】

データ表示装置であって、

複数の文書から構成される文書群中に含まれるキーワードを抽出するキーワード抽出手段と、

前記抽出された各キーワードの、前記文書群中に出現する頻度を算出する頻度算出手段と、

前記各キーワードを前記各キーワードが属するクラスターにクラスタリングするクラスタリング手段と、

前記算出された頻度がより高いキーワードを含む文書の順に、前記文書群を構成する文書の文書データをソートするデータソート手段と、

前記データソート手段によってソートされた文書データに対応付けて、前記クラスターに属するキーワードのうち前記ソートされた文書データを持つ文書が含むキーワードを、

10

20

30

40

50

前記クラスター毎に画面表示する表示手段とを備える
ことを特徴とするデータ表示装置。

【請求項 14】

請求項 13 に記載のデータ表示装置において、
前記データソート手段は、さらに、前記クラスターを前記クラスターに属するキーワードの数について降順又は昇順にソートすることを特徴とするデータ表示装置。

【請求項 15】

データ表示方法であって、
複数の文書から構成される文書群中に含まれるキーワードを抽出し、
前記各キーワードを前記各キーワードが属するクラスターにクラスタリングし、
前記クラスターに属するキーワードのうちの前記文書群を構成する文書に含まれるキーワード、又は、各文書中における前記クラスターに属するキーワードの有無を示す情報を、前記クラスター毎に画面表示することを特徴とするデータ表示方法。

10

【請求項 16】

データ表示プログラムであって、
コンピュータに、
複数の文書から構成される文書群中に含まれるキーワードを抽出する処理と、
前記各キーワードを前記各キーワードが属するクラスターにクラスタリングする処理と

20

、
前記クラスターに属するキーワードのうちの前記文書群を構成する文書に含まれるキーワード、又は、各文書中における前記クラスターに属するキーワードの有無を示す情報を、前記クラスター毎に画面表示する処理とを実行させることを特徴とするデータ表示プログラム。

【発明の詳細な説明】

【技術分野】

【0001】

本発明は、データ表示技術に関し、特に、文書の大凡の内容を一見して把握できるような表示を行うデータ表示装置、データ表示方法及びデータ表示プログラムに関する。

30

【背景技術】

【0002】

従来から、文書中に含まれるキーワードの頻度を算出し、算出された頻度を表示する技術は存在した。

【0003】

また、例えば、データを表形式で表示する技術について、下記の非特許文献 1 に記載されている。

【非特許文献 1】 知りたい操作がすぐわかる 標準 Excel 全機能 Bible 2003, 村田吉徳著, 技術評論社, 2004.2.1 発行

【発明の開示】

40

【発明が解決しようとする課題】

【0004】

しかし、上記従来技術は、文書群中に含まれるキーワードをキーワード群（クラスター）毎にクラスタリングし、当該クラスターに属するキーワードのうちの文書群を構成する文書に含まれるキーワード、又は、各文書中における当該クラスターに属するキーワードの有無を示す情報を、当該各文書と当該クラスターとに対応付けて画面表示することを行っていない。

【0005】

従って、従来技術では、画面表示された各文書データが含むキーワードを参照して、当該文書データを持つ文書の大凡の内容を一見して把握することは困難である。

50

【0006】

本発明は、文書の大凡の内容を一見して把握できるような表示を行うデータ表示装置、データ表示方法及びデータ表示プログラムの提供を目的とする。

【課題を解決するための手段】

【0007】

本発明のデータ表示装置は、複数の文書から構成される文書群中に含まれるキーワードを抽出するキーワード抽出手段と、前記抽出されたキーワードを各キーワードが属するクラスターにクラスタリングするクラスタリング手段と、前記クラスターに属するキーワードのうちの前記文書群を構成する文書に含まれるキーワード、又は、各文書中における前記クラスターに属するキーワードの有無を示す情報を、前記各文書と前記クラスターとに
10 対応付けて画面表示する表示手段とを備える。

【0008】

好ましくは、本発明のデータ表示装置において、前記クラスタリング手段は、前記各キーワードのベクトル空間上の位置を示す位置ベクトルを生成し、生成された前記位置ベクトルが示す前記各キーワードの位置情報に基づいて、前記各キーワードが属するクラスターを決定する。

【0009】

好ましくは、本発明のデータ表示装置において、前記表示手段は、予め定められた順序に基づいて、前記各クラスターを並び替えて画面表示する。

【0010】

好ましくは、本発明のデータ表示装置において、前記表示手段は、各クラスターに属するキーワード又は前記各クラスターに属するキーワードが出現する文書の分布に基づいて、前記各クラスターを並び替えて画面表示する。
20

【0011】

好ましくは、本発明のデータ表示装置において、表示手段は、各クラスターに属するキーワードが出現する文書数について昇順又は降順に前記各クラスターを並び替えて画面表示する。

【0012】

好ましくは、本発明のデータ表示装置において、前記クラスタリング手段は、前記各キーワードの位置ベクトルの平均を前記各キーワードが属するクラスターの位置ベクトルとし、各クラスターの位置ベクトルに基づいて、クラスター同士の距離を求め、求めた距離が近いクラスター同士が近い位置に並び、求めた距離が遠いクラスター同士が離れた位置に並ぶように、前記各クラスターを並び替え、前記表示手段は、前記並び替えられた各クラスターに属するキーワードのうちの前記文書群を構成する文書に含まれるキーワード、又は、各文書中における前記並び替えられた各クラスターに属するキーワードの有無を示す情報を、前記各文書と前記並び替えられた各クラスターとに対応付けて画面表示する。
30

【0013】

好ましくは、本発明のデータ表示装置において、前記クラスタリング手段は、ベクトル空間上に並んで配置された複数の基準の位置ベクトルを予め定義し、前記各クラスターの位置ベクトルを、前記各クラスターの位置ベクトルとの距離が最も近い基準の位置ベクトルの近くに配置することにより、前記各クラスターを並び替える。
40

【0014】

好ましくは、本発明のデータ表示装置において、前記クラスタリング手段は、前記各文書がどのクラスターのどの単語を何個含んでいるかの情報を求め、求めた情報に基づいて、前記各文書のベクトル空間上の位置を示す位置ベクトルを求め、求めた前記各文書の位置ベクトルに基づいて、文書同士の距離を求め、求めた距離が近い文書同士が近い位置に並び、求めた距離が遠い文書同士が離れた位置に並ぶように、各文書を並び替える。

【0015】

好ましくは、本発明のデータ表示装置において、前記表示手段は、並ぶ順序がより上位のクラスターに属するキーワードを含む文書の順に、前記文書に含まれるキーワード、又
50

は、前記文書中における前記クラスターに属するキーワードの有無を示す情報を画面表示する。

【0016】

好ましくは、本発明のデータ表示装置において、前記表示手段は、前記クラスタリング手段によってクラスタリングされたクラスターを選択し、選択されたクラスターに対応付けて、前記クラスターに属するキーワードのうちの前記文書群を構成する文書に含まれるキーワード、又は、各文書中における前記クラスターに属するキーワードの有無を示す情報を画面表示する。

【0017】

好ましくは、本発明のデータ表示装置において、前記クラスタリング手段は、前記各キーワードのベクトル空間上の位置を示す位置ベクトルを生成し、生成された前記位置ベクトルが示す前記各キーワードの位置情報に基づいて、キーワード同士の距離を求め、求めた距離が近いキーワード同士が近い位置に並び、求めた距離が遠いキーワード同士が離れた位置に並ぶように、各キーワードを並び替え、前記表示手段は、並び替えられたキーワードのうちの前記文書群を構成する文書に含まれるキーワード、又は、各文書中における前記並び替えられたキーワードの有無を示す情報を、前記各文書に対応付けて画面表示する。

10

【0018】

好ましくは、本発明のデータ表示装置において、前記表示手段は、各クラスターを画面上に並ばせる順序を選択し、前記選択された順序がより上位のクラスターに属するキーワードを含む文書の順に、前記文書に含まれるキーワード、又は、前記文書中における前記クラスターに属するキーワードの有無を示す情報を画面表示する。

20

【0019】

また、本発明のデータ表示装置は、複数の文書から構成される文書群中に含まれるキーワードを抽出するキーワード抽出手段と、前記抽出された各キーワードの、前記文書群中に出現する頻度を算出する頻度算出手段と、前記各キーワードを前記各キーワードが属するクラスターにクラスタリングするクラスタリング手段と、前記算出された頻度がより高いキーワードを含む文書の順に、前記文書群を構成する文書の文書データをソートするデータソート手段と、前記データソート手段によってソートされた文書データに対応付けて、前記クラスターに属するキーワードのうち前記ソートされた文書データを持つ文書が含むキーワードを、前記クラスター毎に画面表示する表示手段とを備える。

30

【0020】

好ましくは、本発明のデータ表示装置において、前記データソート手段は、さらに、前記クラスターを前記クラスターに属するキーワードの数について降順又は昇順にソートする。

【0021】

また、本発明のデータ表示方法は、複数の文書から構成される文書群中に含まれるキーワードを抽出し、前記各キーワードを前記各キーワードが属するクラスターにクラスタリングし、前記クラスターに属するキーワードのうちの前記文書群を構成する文書に含まれるキーワード、又は、各文書中における前記クラスターに属するキーワードの有無を示す情報を、前記クラスター毎に画面表示する。

40

【0022】

また、本発明のデータ表示プログラムは、コンピュータに、複数の文書から構成される文書群中に含まれるキーワードを抽出する処理と、前記各キーワードを前記各キーワードが属するクラスターにクラスタリングする処理と、前記クラスターに属するキーワードのうちの前記文書群を構成する文書に含まれるキーワード、又は、各文書中における前記クラスターに属するキーワードの有無を示す情報を、前記クラスター毎に画面表示する処理とを実行させる。

【発明の効果】

【0023】

50

本発明のデータ表示装置、データ表示方法及びデータ表示プログラムは、複数の文書から構成される文書群中に含まれるキーワードを抽出し、上記各キーワードを上記各キーワードが属するクラスターにクラスタリングし、当該クラスターに属するキーワードのうちの上記文書群を構成する文書に含まれるキーワード、又は、各文書中における当該クラスターに属するキーワードの有無を示す情報を、当該各文書と当該クラスターとに対応付けて画面表示する。従って、本発明によれば、文書の大凡の内容を一見して把握できるような表示を行うことが可能となる。

【発明を実施するための最良の形態】

【0024】

以下に、図を用いて、本発明の実施の形態について説明する。図1は、本発明のシステム構成の一例を示す図である。データ表示装置1は、文書群中の文書の文書データをソートして画面表示する処理装置である。 10

【0025】

データ表示装置1は、キーワード抽出部11、頻度算出部12、クラスタリング部13、データソート部14、表示部15、クラスターソート指定部16、書誌データデータベース(DB)17を備える。

【0026】

キーワード抽出部11は、書誌データDB17に蓄積されている文書群に含まれるキーワードを抽出する。キーワード抽出部11によるキーワードの抽出手法については、後述する。 20

【0027】

頻度算出部12は、キーワード抽出部11によって抽出された各キーワードの、書誌データDB17に蓄積されている文書群中に出現した頻度を算出する。ここで、キーワードの頻度とは、例えば、キーワードが出現する文書の数を意味する。例えば、キーワード「本塁打記録」を含む文書数が20である場合には、算出されるキーワード「本塁打記録」の頻度は20である。また、本発明の実施の形態においては、文書群でのキーワードの出現回数をキーワードの頻度とする構成を採ることもできる。

【0028】

クラスタリング部13は、キーワード抽出部11によって抽出されたキーワードを、1又は複数のクラスターにクラスタリング(分類)する。ここで、クラスターとは、1又は複数のキーワードで構成されるキーワード群である。クラスタリング部13は、例えば、後述するように、各キーワードの位置(ベクトル空間における位置)を示すベクトル(位置ベクトル)を生成する。そして、クラスタリング部13は、生成されたベクトルが示す各キーワードの位置情報に基づいて、各キーワードが属するクラスターを決定する。例えば、キーワード「リーグ」、「米大リーグ」、「大リーグ」をあるクラスターに属するキーワードとしてクラスタリングし、キーワード「試合」、「チーム」、「スタジアム」を、他のクラスターに属するキーワードとしてクラスタリングする。 30

【0029】

データソート部14は、書誌データDB17に蓄積されている文書から、文書データ(例えば、文書中のテキストデータ、文書のタイトル、著者名等)を抽出し、抽出した文書データをソートする。 40

【0030】

すなわち、データソート部14は、まず、抽出した文書データを図示しないバッファ中に格納する。そして、データソート部14は、頻度算出部12によって算出された各キーワードの頻度がより高いキーワードを多く含む文書の順に、上記バッファ中に格納された文書の文書データをソートする。データソート部14は、後述するクラスターソート指定部16によって指示された、各クラスター(キーワード群)のソート内容に従って、各クラスターをソートするようにしてもよい。

【0031】

なお、本発明の一実施形態によれば、データソート部14は、書誌データDB17から 50

抽出した各文書から各文書に関連する日付（例えば、発行日）のデータを抽出し、文書データ（例えば、文書のテキストデータ、文書のタイトル、著者名等）を日付について降順又は昇順にソートする構成を採ってもよい。

【0032】

表示部15は、データソート部14によってソートされた文書データに対応付けて、上記クラスター中のキーワードのうち当該ソートされた文書データを持つ文書が含むキーワードを、クラスター毎に画面表示する。表示部15は、データソート部14によってソートされた文書データに対応付けて、データソート部14によってソートされたクラスター中のキーワードのうち当該ソートされた文書データを持つ文書が含むキーワードを、データソート部14によってソートされた各クラスター毎に画面表示するようにしてもよい。

10

【0033】

クラスターソート指定部16は、データソート部14による各クラスターのソート内容を指定する。例えば、クラスターソート指定部16は、各クラスターを、クラスター中のキーワードの数について降順にソートすることを指示する制御情報をデータソート部14に対して送出する。クラスターソート指定部16は、各クラスターを、クラスター中のキーワードの数について昇順にソートすることを指示する制御情報をデータソート部14に対して送出するようにしてもよい。また、クラスターソート指定部16は、例えば、後述する各クラスター間の距離に基づいて各クラスターをソートすることを指示する制御情報をデータソート部14に対して送出するようにしてもよい。例えば、クラスターソート指定部16は、後述する各クラスター間の距離が近いクラスター同士が直線上に結合するように各クラスターをソートすることを指示する制御情報をデータソート部14に対して送出する。当該制御情報を受けたデータソート部14は、例えば、距離が最も近いクラスター同士をまず結合し、結合したクラスターによって構成される、複数のクラスター群を作成する。そして、データソート部14は、例えば、あるクラスター群におけるクラスターのいずれかと距離が最も近い他のクラスター群のクラスターを、当該距離が最も近いクラスターと結合する処理を繰り返して、クラスターを直線上に結合する。

20

【0034】

書誌データDB17には、大量の文書（書誌データ）が蓄積されている。

【0035】

本発明の一実施形態によれば、頻度算出部12は、キーワード抽出部11が抽出した各キーワードについて算出した頻度と、各キーワードの文字数とに基づいて、各キーワードのスコアを算出するようにしてもよい。各キーワードのスコアは、例えば、各キーワードの文字数に頻度を乗じた値として算出される。

30

【0036】

本発明の実施の形態においては、頻度算出部12が、キーワード抽出部11によって抽出された各キーワードの文字数を用いずに、各キーワードについて算出された頻度に基づいて、所定の計算式を用いて、各キーワードのスコアを算出する構成を採ってもよい。

【0037】

例えば、頻度算出部12は、以下に示すような、TF/IDF法を用いたスコアの算出方法又はOkapiのウェイト法を用いて、各キーワードのスコアを算出する。

40

【0038】

（TF/IDF法を用いたスコアの算出方法）

一般に、重要なキーワードを含む文書の検索には、主にTF/IDF法が用いられる。ここで、TFとは、一般に、ある文書でのあるキーワードの出現回数を意味し、IDFとは、一般に、予め用意された多数の文書のうち、上記キーワードが出現する文書数の逆数を意味する。

【0039】

一般に、TF/IDF法では、以下の式で算出されるScore(D)が高い文書を検索結果として出力する。

【0040】

50

$$\text{Score}(D) = (tf(w, D) \times \log(N/df(w)))$$

上記の式において、 w は、ユーザが入力するキーワード、 $tf(w, D) \times \log(N/df(w))$ を w W で加算することを意味する。また、 $tf(w, D)$ は、文書 D での w の出現回数であり、 $df(w)$ は、全文書において w が出現した文書の数であり、 N は、文書の総数である。

【0041】

TF/IDF法の本発明への適用に当たっては、例えば、上記文書 D を、書誌データDB17に蓄積されている文書群として、 $tf(w, D)$ を算出する。また、例えば、書誌データDB17とは別のデータベース(図示を省略)に蓄積されている大量の文書群を、上記 $df(w)$ の意味の説明において記述した「全文書」として、 $df(w)$ を算出する。

10

【0042】

そして、算出された $tf(w, D)$ と $\log(N/df(w))$ との積を、各キーワード w のスコアとして算出する。

【0043】

(Okapiのウェイト法を用いたスコアの算出方法)

一般に、Okapiのウェイト法(下記の文献(1)参照)では、以下の式で算出される $\text{Score}(D)$ が高い文書を検索結果として出力する。

【0044】

文献(1):村田真樹,馬青,内元清貴,小作浩美,内山将夫,井佐原均,位置情報と分野情報を用いた情報検索,自然言語処理(言語処理学会誌),2000年4月,7巻,2号,p.141~p.160

20

【0045】

【数1】

$$\text{Score}(D) = \sum_{w \in W} \left[\frac{tf(w, D)}{\frac{\text{length}(D)}{\Delta} + tf(w, D)} \times \log \frac{N}{df(w)} \right]$$

30

【0046】

ここで、 w は、ユーザが入力するキーワード、 W は、ユーザが入力するキーワードの集合を意味する。また、 $tf(w, D)$ は、文書 D での w の出現回数であり、 $df(w)$ は、全文書において w が出現した文書の数であり、 N は、文書の総数である。また、 $\text{length}(D)$ は、文書 D の長さ(文字列単位)である。 Δ は、全文書における文書の長さの平均である。

【0047】

Okapiのウェイト法の本発明への適用に当たっては、例えば、上記文書 D を、書誌データDB17に蓄積されている文書群として、

40

【0048】

【数2】

$$\frac{tf(w, D)}{\frac{\text{length}(D)}{\Delta} + tf(w, D)}$$

【0049】

を算出する。算出された値を tf 項とする。

50

【0050】

また、例えば、書誌データDB17とは別のデータベース(図示を省略)に蓄積されている大量の文書群を、上記df(w)の意味の説明において記述した「全文書」として、 $\log(N/df(w))$ を算出する。算出された $\log(N/df(w))$ をidf項とする。そして、算出されたtf項とidf項との積を、各キーワードwのスコアとして算出する。

【0051】

また、本発明の一実施形態によれば、データソート部14は、頻度算出部12によって算出された各キーワードのスコアがより高いキーワードを多く含む文書の順に、各文書の文書データをソートするようにしてもよい。

10

【0052】

次に、クラスタリング部13による、キーワードのクラスタリング処理について具体的に説明する。クラスタリングには、以下に示すような様々な方法がある。

(階層クラスタリング(ボトムアップクラスタリング)による方法)

距離が最も近い成員同士を結合していき、クラスターを作る。そして、距離が最も近いクラスター同士を結合する。成員とは、クラスタリングの対象となるキーワード(単語)であって、あるクラスターに属するキーワード(単語)である。クラスター間の距離の定義は様々ある。例えば、クラスターAとクラスターBとの距離を、クラスターAの成員(すなわち、クラスターAに属するキーワード)とクラスターBの成員(すなわち、クラスターBに属するキーワード)との距離の中で最も小さいものとしてもよい。ここで、成員と成員との距離とは、ベクトルで表現される後述する成員の位置間の距離である。また、例えば、クラスターAとクラスターBとの距離を、クラスターAの成員とクラスターBの成員との距離の中で最も大きいものとしてもよい。また、例えば、クラスターAとクラスターBとの距離を、全てのクラスターAの成員とクラスターBの成員との距離の平均としてもよい。また、全てのクラスターAの成員の位置の平均をクラスターAの位置とし、全てのクラスターBの成員の位置の平均をクラスターBの位置とし、当該クラスターAの位置とクラスターBの位置との距離をクラスターAとクラスターBとの距離としてもよい。

20

【0053】

(ワード法による方法)

以下に示すWを定義する。

30

$$W = \sum_{i=1}^g \sum_{j=1}^{n_i} (x(i, j) - ave_x(i))^2$$

^は指数を意味する。例えば、上記の式における1つ目の \sum は、 $i=1$ から $i=g$ までの加算、2つ目の \sum は、 $j=1$ から $j=n_i$ までの加算を意味する。また、 $x(i, j)$ は、 i 番目のクラスターの j 番目の成員の位置、 $ave_x(i)$ は、 i 番目のクラスターの全ての成員の位置の平均を意味する。クラスター同士を結合していくと、Wの値が増加するが、ワード法では、Wの値がなるべく大きくならないようにクラスター同士を結合していく。

【0054】

(クラスタリングの終了条件)

予めクラスターの個数を決めておいて、クラスターの個数が当該予め決められた数になったときに、クラスター同士を結合するのをやめるようにしてもよい。また、予め距離の閾値を決めておいて、その閾値数以上離れているクラスター同士を結合するのをやめるようにしてもよい。

40

【0055】

(各成員の位置)

各成員(単語)の位置は、後述するように、各成員に関する種々の情報(例えば、各成員の属性情報)を用いて求める。上記各成員に関する種々の情報に基づいて、ベクトルの次元を決定する。そして、各成員に関する種々の情報に基づいて、上記決定された次元を持つベクトルの要素の値を求めてベクトル(位置ベクトル)を生成する。生成したベクトルは、各成員の位置を示している。各成員(単語)に関する種々の情報としては、例えば

50

、以下に示すものがある。

- ・ 単語に含まれる文字の種類（例えば、ひらがな、カタカナ、漢字、それ以外が、それぞれあるかないか）

- ・ 単語の長さ
- ・ 単語の語義
- ・ 単語の共起語
- ・ 単語の共起データ

(1) 単語に含まれる文字の種類（例えば、ひらがな、カタカナ、漢字、それ以外が、それぞれあるかないか）

例えば、ある成員の文字の種類を表すために、ひらがなのみからなる単語か否か、カタカナのみからなる単語か否か、漢字のみからなる単語か否か、ひらがなのみ、又は、カタカナのみ、又は、漢字のみからなる単語以外の単語か否かという、4次元のベクトルの要素を用意（例えば、設定）し、当該ベクトルの要素に設定される値（例えば、1又は0）によって決まるベクトルを作成する。成員がひらがなのみからなる単語に該当すれば1の値を、該当しなければ0の値を対応するベクトルの要素に設定する。また、成員がカタカナのみからなる単語に該当すれば1の値を、該当しなければ0の値を対応するベクトルの要素に設定する。また、成員が漢字のみからなる単語に該当すれば1の値を、該当しなければ0の値を対応するベクトルの要素に設定する。また、成員がひらがなのみ、又は、カタカナのみ、又は、漢字のみからなる単語以外の単語に該当すれば1の値を、該当しなければ0の値を対応するベクトルの要素に設定する。

10

20

(2) 単語の長さ

例えば、1次元のベクトルの要素を用意し、成員（単語）の文字の個数をカウントし、カウントされた成員（単語）の文字の個数を当該ベクトルの要素に設定して、単語の長さの情報を表現してもよい。

(3) 単語の語義

例えば、予め記憶手段内に記憶された、以下のような分類語彙表を用意する。

あ, あ, 4.310, 1, 10, *,
 あ, 亜, 1.104, 2, 40,,
 あ, 亜, 3.100, 10, 40,,
 ああ, ああ, 3.100, 3, 40, *,
 ああ, ああ, 4.310, 1, 20, *,
 ああくとう, アーク燈, 1.460, 2, 70,,
 ああす, アース, 1.462, 6, 10,,
 ああち, アーチ, 1.442, 2, 20,,
 ああむほおる, アームホール, 1.184, 5, 30,,
 あある, アール, 1.1961, 4, 10,,
 あい, 愛, 1.3020, 9, 10, *,
 あい, 相, 3.112, 1, 10, *,
 あい, 藍, 1.502, 6, 40,,
 あいいく, 愛育, 1.3642, 1, 40,,
 あいいん, 愛飲, 1.3332, 3, 60,,
 あいいん, 合印, 1.3114, 1, 30, Y,
 あいうち, あい打ち, 1.357, 4, 30,,
 あいかぎ, 合鍵, 1.454, 8, 50,,
 あいかわらず, 相変らず, 3.165, 2, 10, *,
 あいかん, 哀歓, 1.3011, 4, 60,,
 あいがん, 哀願, 1.366, 1, 100,,
 あいがん, 愛断, 1.3852, 2, 10,,
 あいぎ, 合着, 1.421, 4, 40,,
 あいきょう, 愛郷, 1.3020, 11, 170,,

30

40

50

あいきょう, 愛嬌, 1.3030, 4, 40, ,

上記の", "で区切ってある情報は、それぞれ、単語の読み、単語の見出し語、単語の分類番号、単語の分類番号の下位番号1、単語の分類番号の下位番号2、標本使用頻度が7以上の単語かどうかを示す情報である。

【0056】

上記の分類語彙表中の", "で区切ってある情報を各桁とし、例えば、上位3桁を意味分類と仮定して、その上位3桁の種類の数だけベクトルの次元を用意する。そして、各成員について、当該意味分類と合致したベクトルの要素の値を1に、それ以外を0に設定することによって、各成員の位置を示すベクトルを作成する。上述した分類語彙表中の情報を利用して各成員のベクトルを作成する方法では、上位3桁を意味分類としたが、他の桁を意味分類にしてもよいし、当該上位3桁と上記他の桁とを合わせた複数の桁の種類の数だけベクトルの次元を用意し、各成員について、当該意味分類と合致したベクトルの要素の値を1に、それ以外を0に設定することによって、各成員の位置を示すベクトルを作成するようにしてもよい。

10

【0057】

各成員(単語)に関する種々の情報として、単語に含まれる文字の種類、単語の長さ、単語の語義を説明したが、本発明の一実施例によれば、各成員に関する他の情報も追加で利用して、各成員の位置を示すベクトルを生成してもよい。また、逆に、上述した情報(単語に含まれる文字の種類、単語の長さ、単語の語彙)の全てを用いるのではなく、それらの情報の一部を用いてベクトルを作成するようにしてもよい。例えば、単語の長さの情報を用いずに、ベクトルを作成するようにしてもよい。また、例えば、単語に含まれる文字の種類と単語の語義のみを用いて、ベクトルを作成するようにしてもよい。

20

(4) 単語の共起語

成員(単語)の共起語を求めて、単語の種類の情報に基づいて決まるベクトルの次元を用意する。そして、当該単語の共起語に合致したベクトルの要素に1を、当該単語の共起語に合致しないベクトルの要素に0を設定する。単語の共起語としては、例えば、図1に示す書誌データDB17中に格納されている書誌データのうち、同じ書誌データ中に成員と共起して出現する単語を当該成員の共起語とする。本発明の一実施形態によれば、例えば、書誌データDB17中に格納されている書誌データとは異なるデータを用意し、当該データにおいて、成員が出現する領域と同一の領域(例えば、同一文、又は、同一段落、又は、同一データレコード等)に出現する単語を、当該成員の共起語とするようにしてもよい。

30

【0058】

また、上述した単語の共起語を用いてベクトルを作成する方法では、共起語として出現していれば、ベクトルの要素に1を、共起語として出現していない場合に、ベクトルの要素に0を設定するが、ベクトルの要素には、共起語として出現した回数の値を設定するようにしてもよい。

(5) 単語の共起データ

例えば、書誌データ17DB中に格納された書誌データの行数だけのベクトルの次元を用意し、成員(単語)がその書誌データのある行に出現した場合、当該行に対応するベクトルの要素に1を、出現しない場合に、当該行に対応するベクトルの要素に0を入れる。

40

【0059】

単語の語義、単語の共起語、単語の共起データについては、意味分類、共起語、共起データの数が多いため、本発明の一実施形態によれば、既存のLSI(Latent semantic index)などの次元圧縮の技術を使って、ベクトルの次元を減らすようにしてもよい。

【0060】

次にトップダウンのクラスタリング(非階層クラスタリング)の方法を説明する。

(最大距離アルゴリズムによるクラスタリング)

ある成員と、当該成員と距離が最も離れた成員を求め、これらの成員をそれぞれのクラ

50

スターの中心とする。次に、それぞれのクラスターの中心と各成員との距離の最小値を各成員の距離とし、その距離が最も大きい成員を新たなクラスターの中心とする。当該クラスターの中心を求める処理を繰り返す。例えば、予め定めた数のクラスターになったときに、当該クラスターの中心を求める処理の繰り返しをやめる。また、例えば、クラスター間の距離が予め定めた数以下になったときに、当該クラスターの中心を求める処理の繰り返しをやめる。

【0061】

また、クラスターの良さを例えばAIC情報量基準などで評価して、評価によって求めた値と予め定めた閾値との比較結果に基づいて、当該クラスターの中心を求める処理の繰り返しをやめるようにしてもよい。上記の最大距離アルゴリズムによるクラスタリングによれば、各成員は、各成員と最も近いクラスター中心を持つクラスターの成員となる。
(k平均法)

10

例えば、以下に示すk平均法によって、予め定めた個数(k個)にクラスタリングする。まず、k個の成員をランダムに選択し、選択されたk個の成員をクラスターの中心とする。そして、各成員を、当該各成員に最も近いクラスター中心を持つクラスターの成員とする。

【0062】

次に、クラスター内の各成員の平均の位置に最も近い成員を、それぞれのクラスターの中心とする。そして、各成員を、当該各成員に最も近いクラスター中心を持つクラスターの成員とする。また、クラスター内の各成員の平均の位置に最も近い成員をそれぞれのクラスターの中心とする。上記のクラスターの中心を求める処理を繰り返し、クラスターの中心が移動しなくなったときに、クラスターの中心を求める処理の繰り返しをやめる。本発明の一実施形態によれば、予め定めた回数だけクラスターの中心を求める処理を繰り返してやめるようにしてもよい。そして、最終的なクラスター中心を持つクラスターを決定する。そして、各成員を、当該各成員が最も近いクラスター中心を持つクラスターの成員とする。上記の手法によって、成員のクラスタリングをする。本発明において用いるクラスタリングの方法は、上述した方法に限定されるものではない。本発明に係るデータ表示装置1は、上述したクラスタリングの方法以外の様々な方法を用いて、キーワードのクラスタリングをするようにしてもよい。例えば、本発明の一実施形態によれば、例えば、距離の近い成員を直線上に順に結合する。すなわち、距離が最も近い成員同士をまず結合し、結合した成員によって構成される、複数の成員群を作成する。ある成員群における成員のいずれかと距離が最も近い他の成員群の成員を、当該距離が最も近い成員と結合する処理を繰り返して、成員を直線上に結合する。そして、直線上に結合された成員からなるリストを、予め定めた数の成員毎に区切って、各成員のクラスタリングをするようにしてもよい。

20

30

【0063】

また、本発明の一実施形態によれば、予めデータ表示装置1内の記憶手段(図示を省略)内に、単語と単語が属する分類(クラスター)との対応情報を予め記憶させておき、クラスタリング部13が、当該記憶手段内の、単語と単語が属する分類(クラスター)との対応情報に基づいて、各単語をクラスタリングするようにしてもよい。

40

【0064】

次に、データソート部14による、抽出した文書の文書データのソート処理について、具体的に説明する。上述したように、データソート部14は、例えば、頻度算出部12によって算出された頻度がより高いキーワードを含む文書の順に、各文書の文書データをソートする。

【0065】

例えば、文書Aが、頻度が最も高いキーワード「本塁打記録」と、頻度が2番目に高いキーワード「本塁打」と、頻度が3番目に高いキーワード「リーグ」とを含んでいるものとし、また、例えば、文書Bが、頻度が最も高いキーワード「本塁打記録」と頻度が3番目に高いキーワード「リーグ」と、頻度が4番目に高いキーワード「マグワイア」とを含

50

んでいるものとする。文書 A は、文書 B に含まれない、頻度が 2 番目に高いキーワードを含んでいる。この場合、文書 A は、文書 B に比べて、頻度がより高いキーワードを含んでいる。

【0066】

頻度がより高いキーワードを含んでいるということ、さらに具体的に説明する。例えば、各キーワードを頻度について降順に並べ、文書があるキーワードを含む場合に、そのキーワードにビット論理「1」を割り当て、文書があるキーワードを含まない場合に、そのキーワードにビット論理「0」を割り当てる。そして、各キーワードに割り当てられたビット論理によって構成される 2 進数を求める。

【0067】

例えば、「本塁打記録」、「本塁打」、「リーグ」、「マグワイア」、・・・の順にキーワードが並ぶとすると、上記の文書 A について求められる 2 進数は、「1110・・・」であり、文書 B について求められる 2 進数「1011・・・」より大きな数となる。

【0068】

ある文書が、頻度がより高いキーワードを含んでいるということは、上記のように、例えば、頻度について降順に並んだ各キーワードを 2 進数の各桁とし、文書に含まれるキーワードにビット論理「1」を、文書に含まれないキーワードにビット論理「0」を割り当てた場合に構成される 2 進数が、より大きい数であることを意味している。

【0069】

以下に、キーワード抽出部 11 によるキーワードの抽出方法について説明する。

(1) 形態素解析を用いた単語の認識による手法

まず、キーワード抽出部 11 は、書誌データ DB 17 に蓄積されている文書について、形態素解析を行い、単語の認識を行う。そして、特定の名詞の単語をキーワードとして取り出す。例えば、名詞だけをキーワードとして取り出す。但し、「こと」、「もの」などの一般的な名詞は、予め収集しておき、それらの名詞がキーワードとしては取り出されないようにしておく。キーワードとしては、名詞だけでなく、動詞などの他の品詞も取り出すこととしてもよい。

【0070】

形態素解析には、例えば、奈良先端大で開発されている形態素解析システムである ChaSen (下記の文献(2)参照)を用いる。

【0071】

文献(2)：形態素解析システム茶筌 (<http://chasen.aist-nara.ac.jp/index.html>.ja)

ChaSen は、日本語文を分割し、さらに、各単語の品詞も推定してくれる。

【0072】

例えば、「学校へ行く」を入力すると、以下の結果を得る。

【0073】

| | | | | | |
|----|------|----|------|-------|----------|
| 学校 | ガッコウ | 学校 | 名詞 - | 一般 | |
| へ | へ | へ | 助詞 - | 格助詞 - | 一般 |
| 行く | イク | 行く | 動詞 - | 自立 | 五段・力行促音便 |
| | EOS | | | | 基本形 |

このように、各行に一個の単語が入るように分割され、各単語に読みや品詞の情報が付与される。

【0074】

また、英語の品詞タグ付けシステムとしては、Brill (下記の文献(3)参照)のものがある。このシステムを用いれば、英語文の各単語の品詞を推定することができる。

【0075】

文献(3)：Eric Brill, Transformation-Based Error-Driven Learning and Natural Language Processing: A Case Study in Part-of-Speech Tagging, Computational Ling

10

20

30

40

50

uistics, Vol. 21, No. 4, p.543-565, 1995.

(2) TF / IDF法などを利用した方法

書誌データDB17に蓄積されている文書について、形態素解析を行い、例えば、名詞だけを取り出す。そして、取り出された各名詞について、前述したTF / IDF法に基づいて算出される所定のスコアを求め、求めたスコアが所定の値よりも大きいものか、スコアが所定の値よりも大きいものから順に所定の値の個数だけ取り出したものをキーワードとする。なお、上記のスコアは、前述したOkapiのウェイト法を用いて算出されるスコアを用いてもよい。

(3) 高精度な既存のキーワード抽出のツールを利用する方法

一般に文書中では複数の単語の組み合わせで複雑な概念を表す場合が多く、文書の内容が専門的な事項に特化すれば、その傾向はさらに顕著なものとなる。そこで、例えば、(a)形態素解析プログラムによる単語分割、(b)複合語の作成、(c)文書中における重要度の計算、という3つのステップを踏むことで、複合語により複雑な概念を表すことが多い専門用語をキーワードとして文書中から抽出することができる。

【0076】

例えば、下記の文献(4)に記載されている手法は、文書から取り出した単名詞について、単名詞の左右に接続する単語の種類数あるいは頻度を用いたスコアを算出し、これら左右のスコアを組み合わせ、単名詞のスコアを算出する。単名詞のスコアに基づいて、単名詞から生成される複合名詞のスコアを算出する。そして、算出された複合名詞のスコアが所定の値より大きいものを、キーワードとして取り出す。本発明においても、文献(4)に記載された手法を用いて、キーワードを抽出する構成を採ることができる。

【0077】

文献(4)：中川裕志、森辰則、湯本紘彰：“出現頻度と接続頻度に基づく専門用語抽出”，自然言語処理、Vol.10 No.1, pp. 27 - 45, 2003年1月

なお、本発明の実施の形態において、キーワード抽出部11によるキーワードの抽出方法は、上述した3つの方法に限定されるものではない。キーワード抽出部11は、他の任意のキーワードの抽出方法を用いてキーワードを抽出することができる。

【0078】

本発明の一実施形態によれば、キーワード抽出部11が、以下に述べる固有表現抽出技術を用いてキーワードを抽出してもよい。固有表現とは、人名、地名、組織名などの固有名詞、金額などの数値表現といった、特定の事物・数量を意味する言語表現のことで、固有表現抽出とは、そういった固有表現を文章中から計算機で自動で抽出する技術である。例えば、「日本の首相は小泉純一郎である」という文に対して固有表現抽出を行なうと、固有表現の「日本」と「小泉純一郎」が地名、人名として、抽出される。キーワード抽出部11は、抽出された固有表現をキーワードとして出力する。

【0079】

そして、本発明の一実施形態によれば、例えば、予めデータ表示装置1内の記憶手段(図示を省略)内に、キーワード(地名、人名等の固有表現)と当該キーワードが属する分類(クラスター)との対応情報を予め記憶させておき、クラスタリング部13が、キーワード抽出部11によって抽出されたキーワード(固有表現)と、当該記憶手段内の、キーワードと当該キーワードが属する分類(クラスター)との対応情報に基づいて、各キーワードをクラスタリングするようにしてもよい。

【0080】

以下に、固有表現抽出の一般的な手法の例について説明する。

(1) 機械学習を用いる手法

機械学習を用いて固有表現を抽出する手法がある(例えば、以下の参考文献(1)参照)。

【0081】

参考文献(1)：浅原正幸，松本裕治，日本語固有表現抽出における冗長的な形態素解析の利用情報処理学会自然言語処理研究会 NL153-7 2002

10

20

30

40

50

まず、例えば、「日本の首相は小泉さんです。」という文を、各文字に分割し、分割した文字について、以下のように、B - LOCATION、I - LOCATION等の正解タグを付与することによって、正解を設定する。以下の一列目は、分割された各文字であり、各文字の正解タグは二列目である。

日 B - LOCATION
 本 I - LOCATION
 の 0
 首 0
 相 0
 は 0
 小 B - PERSON
 泉 I - PERSON
 さ 0
 ん 0
 で 0
 す 0
 。 0

10

上記において、B - ???は、ハイフン以下の固有表現の種類が始まりを意味するタグである。例えば、B - LOCATIONは、地名という固有表現の始まりを意味しており、B - PERSONは、人名という固有表現の始まりを意味している。また、I - ???は、ハイフン以下の固有表現の種類が始まり以外を意味するタグであり、0はこれら以外である。従って、例えば、文字「日」は、地名という固有表現の始まりに該当する文字であり、文字「本」までが地名という固有表現である。

20

【0082】

このように、各文字の正解を設定しておき、このようなデータから学習し、新しいデータでこの正解を推定し、この正解のタグから、各固有表現の始まりと、どこまでがその固有表現かを認識して、固有表現を推定する。

【0083】

この各文字に設定された正解のデータから学習するときには、システムによってさまざまな情報を素性という形で利用する。例えば、

30

日 B - LOCATION

の部分、

日本 - B 名詞 - B

などの情報を用いる。日本 - B は、日本という単語の先頭を意味し、名詞 - B は、名詞の先頭を意味する。単語や品詞の認定には、例えば前述したChaSenによる形態素解析を用いる。ChaSenを用いれば、入力された日本語を単語に分割することができる。例えば、ChaSenは、前述したように、日本語文を分割し、さらに、各単語の品詞も推定してくれる。例えば、「学校へ行く」を入力すると以下の結果を得ることができる。

【0084】

| | | | | | |
|----|------|----|---------------|----------|-----|
| 学校 | ガッコウ | 学校 | 名詞 - 一般 | | |
| へ | へ | へ | 助詞 - 格助詞 - 一般 | | |
| 行く | イク | 行く | 動詞 - 自立 | 五段・力行促音便 | 基本形 |

E O S

40

このように各行に一個の単語が入るように分割され、各単語に読みや品詞の情報が付与される。

【0085】

なお、例えば、上記の参考文献(1)では、素性として、入力文を構成する文字の、文字自体(例えば、「小」という文字)、字種(例えば、ひらがなやカタカナ等)、品詞情報、タグ情報(例えば、「B - PERSON」等)を利用している。

【0086】

50

これら素性を利用して学習する。タグを推定する文字やその周辺の文字にどのような素性
 が出現するかを調べ、どのような素性が出現しているときにどのようなタグになりやすいかを
 学習し、その学習結果を利用して新しいデータでのタグの推定を行なう。機械学習には、
 例えばサポートベクトルマシンを用いる。

【0087】

固有表現抽出には、上記の手法の他にも種々の手法がある。例えば、最大エントロピー
 モデルと書き換え規則を用いて固有表現を抽出する手法がある（参考文献（2）参照）。

【0088】

参考文献（2）：内元清貴，馬青，村田真樹，小作浩美，内山将夫，井佐原均，最大エ
 ントロピーモデルと書き換え規則に基づく固有表現抽出，言語処理学会誌，Vol.7，No.2，
 2000

10

また、例えば、以下の参考文献（3）に、サポートベクトルマシンを用いて日本語固有
 表現抽出を行う手法について記載されている。

【0089】

参考文献（3）：山田寛康，工藤拓，松本裕治，Support Vector Machineを用いた日本
 語固有表現抽出，情報処理学会論文誌，Vol.43，No.1"，2002

（2）作成したルールを用いる手法

人手でルールを作って固有表現を取り出すという方法もある。

【0090】

例えば、

名詞＋「さん」だと人名とする

名詞＋「首相」だと人名とする

名詞＋「町」だと地名とする

名詞＋「市」だと地名とする

などである。

【0091】

本発明の一実施形態によれば、データ表示装置1は、キーワード抽出部11を用いない
 構成を採ることもできる。例えば、文書と文書に含まれるキーワードとが対応付けられた
 データを所定のデータベース（図1では図示を省略）内に蓄積しておき、上記データベ
 ース内に蓄積されているデータから、頻度算出部12が、各キーワードの、上記データベ
 ース中の文書群中に出現した頻度を算出する構成を採ることもできる。

20

30

【0092】

本発明の一実施形態によれば、キーワード抽出部11が、書誌データDB17に蓄積さ
 れている文書群から、単位表現を抽出し、抽出された単位表現が上記文書群において出現
 する箇所を特定し、特定された箇所において、当該単位表現と隣接して出現している数値
 を抽出する。そして、キーワード抽出部11は、抽出された数値と当該単位表現とから構
 成される数値情報をキーワードとして抽出する。例えば、40号等といった数値情報がキ
 ーワードとして抽出される。そして、クラスタリング部13が、抽出された数値情報をク
 ラスタリングするようにしてもよい。

【0093】

本発明の一実施形態によれば、データ表示装置1は、頻度算出部12、データソート部
 14、クラスタソート指定部16が省略された構成を採ってもよい。例えば、キーワー
 ド抽出部11が、書誌データDB17中の文書群中に含まれるキーワードを抽出し、クラ
 スタリング部13が、抽出されたキーワードを各キーワードが属するクラスターにクラ
 スタリングし、表示部15が、当該クラスターに属するキーワードのうちの上記文書群を構
 成する文書に含まれるキーワード、又は、各文書中における当該クラスターに属するキ
 ーワードの有無を示す情報（例えば、符号）を、当該各文書と当該クラスターとに対応付
 けて画面表示するようにしてもよい。また、データ表示装置1が上記の構成を採る場合
 において、クラスタリング部13が、各キーワードのベクトル空間上の位置を示す位置ベクト
 ルを生成し、生成された当該位置ベクトルが示す各キーワードの位置情報に基づいて、各

40

50

キーワードが属するクラスターを決定するようにしてもよい。

【0094】

また、データ表示装置1が上記の構成を採る場合において、表示部15は、予め定められた順序に基づいて、各クラスターを並び替えて画面表示するようにしてもよい。また、表示部15が、各クラスターに属するキーワード又は各クラスターに属するキーワードが出現する文書の分布に基づいて、各クラスターを画面表示するようにしてもよい。ここで、各クラスターに属するキーワードの分布とは、例えば、各クラスターに属するキーワードが出現する文書数の分布、各クラスターに属するキーワード数の分布、又は、各クラスターに属するキーワードが出現する文書数を各クラスターに属する全てのキーワードについて合計した数の分布である。また、表示部15は、各クラスターに属するキーワードが出現する文書数について昇順又は降順に各クラスターを並び替えて画面表示するようにしてもよい。

10

【0095】

また、データ表示装置1において、クラスタリング部13が、各キーワードの位置ベクトルの平均を各キーワードが属するクラスターの位置ベクトルとし、各クラスターの位置ベクトルに基づいて、クラスター同士の距離を求め、求めた距離が近いクラスター同士が近い位置に並び、求めた距離が遠いクラスター同士が離れた位置に並ぶように、各クラスターを並び替え、表示部15が、並び替えられた各クラスターに属するキーワードのうち、の書誌データDB17中の文書群を構成する文書に含まれるキーワード、又は、各文書中における上記並び替えられた各クラスターに属するキーワードの有無を示す情報を、各文書と当該並び替えられた各クラスターとに対応付けて画面表示するようにしてもよい。

20

【0096】

本発明の一実施形態によれば、クラスタリング部13は、例えば、以下に述べる手法によって、各クラスターを並び替えるようにしてもよい。すなわち、クラスタリング部13は、距離が最も近いクラスター同士を結合し、結合したクラスターのリストの端のクラスターのいずれかと距離が最も近いクラスターを、当該リストの端のクラスターと結合する。そして、クラスタリング部13は、さらに、当該クラスターの結合によって得られるリストの端のクラスターのいずれかと距離が最も近いクラスターを結合する。クラスタリング部13は、上記の処理を繰り返すことによって得られるリストにおけるクラスターの並び順に、各クラスターを並び替える。

30

【0097】

また、クラスタリング部13が、ベクトル空間上に並んで配置された複数の基準の位置ベクトルを予め定義し、各クラスターの位置ベクトルを、各クラスターの位置ベクトルとの距離が最も近い基準の位置ベクトルの近くに配置することにより、各クラスターを並び替えるようにしてもよい。すなわち、クラスタリング部13は、例えば基準ベクトルA、B、Cの順に並んで配置された3つの基準ベクトルを定義し、あるクラスターの位置ベクトルとの距離が最も近い基準ベクトルが基準ベクトルAである場合、当該クラスターの位置ベクトルを基準ベクトルAの近くに配置する。同様にして、他のクラスターの位置ベクトルを、距離が最も近い基準ベクトルの近くに配置する。そして、クラスタリング部13は、各クラスターの位置ベクトルが上記配置後の位置ベクトルとなるように各クラスターを並び替える。

40

【0098】

また、データ表示装置1において、クラスタリング部13は、各文書がどのクラスターのどの単語を何個含んでいるかの情報を求め、求めた情報に基づいて、各文書のベクトル空間上の位置を示す位置ベクトルを求め、求めた各文書の位置ベクトルに基づいて、文書同士の距離を求め、求めた距離が近い文書同士が近い位置に並び、求めた距離が遠い文書同士が離れた位置に並ぶように、各文書を並び替えるようにしてもよい。例えば、クラスタリング部13は、各文書がどのクラスターのどの単語を何個含んでいるかの情報に基づいて、当該各文書の位置ベクトルの次元数と、当該位置ベクトルのベクトル要素の値を求めることによって、各文書の位置ベクトルを求める。

50

【0099】

本発明の一実施形態によれば、クラスタリング部13は、例えば、以下に述べる手法によって、各文書を並び替えるようにしてもよい。すなわち、クラスタリング部13は、距離が最も近い文書同士を結合し、結合した文書のリストの端の文書のいずれかと距離が最も近い文書を、当該リストの端の文書と結合する。そして、クラスタリング部13は、さらに、当該文書の結合によって得られるリストの端の文書のいずれかと距離が最も近い文書を結合する。クラスタリング部13は、上記の処理を繰り返すことによって得られるリストにおける文書の並び順に、各文書を並び替える。

【0100】

また、表示部15は、並び順序がより上位のクラスターに属するキーワードを含む文書の順に、当該文書に含まれるキーワード、又は、当該文書中における上記クラスターに属するキーワードの有無を示す情報（例えば、符号）を画面表示するようにしてもよい。例えば、クラスターA、クラスターB、クラスターCの順にクラスターが並び場合、文書AがクラスターAに属するキーワードとクラスターBに属するキーワードとクラスターCに属するキーワードを含み、文書BがクラスターAに属するキーワードとクラスターBに属するキーワードとを含むがクラスターCに属するキーワードを含まないとき、表示部15は、文書Aに含まれるキーワード、文書Bに含まれるキーワードの順に画面表示する。

【0101】

また、表示部15は、クラスタリング部13によってクラスタリングされたクラスターを、ユーザの指定入力に従って選択し、選択されたクラスターに対応付けて、当該クラスターに属するキーワードのうちの、文書群を構成する文書に含まれるキーワード、又は、各文書中における当該クラスターに属するキーワードの有無を示す情報を画面表示するようにしてもよい。すなわち、表示部15は、例えば、選択されないクラスターについては、当該クラスターに属するキーワードのうちの、文書群を構成する文書に含まれるキーワード、又は、各文書中における当該クラスターに属するキーワードの有無を示す情報を画面表示しないようにしてもよい。

【0102】

また、データ表示装置1において、クラスタリング部13は、各キーワードのベクトル空間上の位置を示す位置ベクトルを生成し、生成された当該位置ベクトルが示す各キーワードの位置情報に基づいて、キーワード同士の距離を求め、求めた距離が近いキーワード同士が近い位置に並び、求めた距離が遠いキーワード同士が離れた位置に並びように、各キーワードを並び替え、表示部15は、並び替えられたキーワードのうちの当該文書群を構成する文書に含まれるキーワード、又は、各文書中における当該並び替えられたキーワードの有無を示す情報を、各文書に対応付けて画面表示するようにしてもよい。

【0103】

本発明の一実施形態によれば、クラスタリング部13は、例えば、以下に述べる手法によって、各キーワードを並び替えるようにしてもよい。すなわち、クラスタリング部13は、距離が最も近いキーワード同士を結合し、結合したキーワードのリストの端のキーワードのいずれかと距離が最も近いキーワードを、当該リストの端のキーワードと結合する。そして、クラスタリング部13は、さらに、当該キーワードの結合によって得られるリストの端のキーワードのいずれかと距離が最も近いキーワードを結合する。クラスタリング部13は、上記の処理を繰り返すことによって得られるリストにおけるキーワードの並び順に、各キーワードを並び替える。

【0104】

また、表示部15は、各クラスターを画面上に並ばせる順序を、ユーザの指定入力に従って選択し、選択された順序がより上位のクラスターに属するキーワードを含む文書の順に、文書に含まれるキーワード、又は、文書中における当該クラスターに属するキーワードの有無を示す情報を画面表示するようにしてもよい。

【0105】

図2は、データ表示処理フローの一例を示す図である。まず、キーワード抽出部11が

、書誌データDB17に蓄積されている文書群に含まれるキーワードを抽出する(ステップS1)。例えば、キーワード「本塁打記録」、「本塁打」、「リーグ」、「マグワイア」、「カージナルス」、「米大リーグ」、「大リーグ」、「試合」、「ロジャー」、「単独トップ」、・・・といったキーワードを抽出する。

【0106】

次に、頻度算出部12が、キーワード抽出部11によって抽出された各キーワードの、書誌データDB17に蓄積されている文書群中に出現した頻度を算出する(ステップS2)。

【0107】

例えば、図3の表に示すように、算出されるキーワード「本塁打記録」の頻度は20、キーワード「本塁打」の頻度は15、キーワード「リーグ」の頻度は12、キーワード「マグワイア」の頻度は10、キーワード「カージナルス」の頻度は9、キーワード「米大リーグ」の頻度は8、キーワード「大リーグ」の頻度は7、キーワード「試合」の頻度は6、キーワード「ロジャー」の頻度は5、キーワード「単独トップ」の頻度は4である。なお、図3中に示す頻度は、各キーワードが出現する文書の数である。また、図3中では、頻度が4であるキーワードまでしか示していないが、本発明の実施の形態では、ステップS2において、例えば、頻度3、頻度2、頻度1についても算出され得る。

10

【0108】

次に、データソート部14が、書誌データDB17に蓄積されている各文書の文書データを抽出し、バッファ中に格納する(ステップS3)。例えば、文書データとして、文書のテキストデータが抽出され、バッファ中に格納される。

20

【0109】

また、データソート部14が、頻度算出部12によって算出された各キーワードの頻度がより高いキーワードを多く含む文書の順に、バッファ中に格納された文書の文書データをソートする(ステップS4)。ステップS4においては、例えば、データソート部14は、各キーワードをステップS2において算出された頻度について降順又は昇順にソートし、ソートした各キーワードの情報をバッファ中に保持するようにしてもよい。ステップS4の処理によって、例えば、各文書と各文書に含まれるキーワードとの対応情報がバッファ中に保持される。

【0110】

図4は、データソート部14によってバッファ中に保持される、各文書と各文書に含まれるキーワードとの対応情報の例を示す図である。図4中において、矩形の枠で囲ったキーワードは、データソート部14によって、例えば上記頻度について降順にソートされたキーワードである。

30

【0111】

例えば、論文名が「A」である文書は、図3中に示される頻度が最も高いキーワード「本塁打記録」と、頻度が2番目に高いキーワード「本塁打」と、頻度が3番目に高いキーワード「リーグ」と、頻度が4番目に高いキーワード「マグワイア」と、頻度が5番目に高いキーワード「カージナルス」と、頻度が6番目に高いキーワード「米大リーグ」とを含んでいるものとする。

40

【0112】

また、例えば、論文名が「B」である文書は、頻度が最も高いキーワード「本塁打記録」と、頻度が2番目に高いキーワード「本塁打」と、頻度が3番目に高いキーワード「リーグ」と、頻度が4番目に高いキーワード「マグワイア」と、頻度が5番目に高いキーワード「カージナルス」とを含んでいるが、頻度が6番目に高いキーワード「米大リーグ」を含んでいないものとする。

【0113】

また、例えば、論文名が「C」である文書は、頻度が最も高いキーワード「本塁打記録」と、頻度が2番目に高いキーワード「本塁打」と、頻度が3番目に高いキーワード「リーグ」と、頻度が4番目に高いキーワード「マグワイア」とを含んでいるが、頻度が5番

50

目に高いキーワード「カージナルス」を含んでいないものとする。

【0114】

上記ステップS4において、データソート部14は、例えば、論文名が「A」という文書のテキストデータ、論文名が「B」という文書のテキストデータ、論文名が「C」という文書のテキストデータの順に並ぶようにソートする。ステップS4における処理の結果、例えば図4中に示すような順番で、各文書のテキストデータがソートされる。なお、図4中の、各論文名が記述された行と同一の行に記述されたキーワードは、各論文名を持つ文書が含むキーワードを示している。また、図4中に示す日付は、各文書の発行日を示す情報である。

【0115】

次に、クラスタリング部13が、キーワード抽出部11によって抽出された各キーワードをクラスター毎にクラスタリングする(ステップS5)。例えば、クラスタリング部13は、図4中の矩形の枠で囲った各キーワードに含まれるキーワードのうち、キーワード「本塁打記録」、「本塁打」、「年間最多本塁打記録」、「年間ホームラン数」、「年間ホームラン」、「号本塁打」、「ホームラン」、「46本」、「50本」、「54本」、「47号」、「50号」、「51号」、「52号」、「54号」(本塁打記録、本塁打以外のキーワードは、図示を省略)を、クラスターAとしてクラスタリングする。また、クラスタリング部13は、例えば、キーワード「リーグ」、「米大リーグ」、「大リーグ」を、クラスターBとしてクラスタリングする。また、クラスタリング部13は、例えば、キーワード「マグワイア」、「マグワイア一塁手」、「マーク」、「ロジャー」、「ソーサ」、「ソーサ外野手」、「サミー」(マグワイア一塁手、マーク、ソーサ、ソーサ外野手、サミーについては、図示を省略)を、クラスターCとしてクラスタリングする。また、クラスタリング部13は、例えば、キーワード「カージナルス」、「ヤンキース」(ヤンキースについては、図示を省略)を、クラスターDとしてクラスタリングする。また、クラスタリング部13は、例えば、キーワード「試合」、「チーム」、「スタジアム」(チーム、スタジアムについては、図示を省略)を、クラスターEとしてクラスタリングする。また、クラスタリング部13は、例えば、キーワード「単独トップ」、「トップ」、「9位タイ」(トップ、9位タイについては、図示を省略)を、クラスターFとしてクラスタリングする。

【0116】

そして、表示部15が、データソート部14によってソートされた文書データに対応付けて、上記クラスター中のキーワードのうち当該ソートされた文書データを持つ文書が含むキーワードを、クラスター毎に画面表示する(ステップS6)。

【0117】

例えば、図5に示す表示画面に示すように、表示部15は、各文書(図5では論文名が「A」である文書～論文名が「E」である文書)のテキストデータと、太線の矩形の枠で囲ったクラスター(図5では、クラスターA～クラスターF)とを画面表示する。また、図5に示す表示画面が、例えばセルで構成されている場合を想定すると、表示部15は、例えば、論文名が「A」である文書のテキストデータ(「米大リーグ、カージナルスのマーク・マグワイア一塁手が・・・」)が表示されるセルに対応する行とクラスターAが表示されるセルに対応する列とが交差するセルの位置に、当該クラスターA中のキーワードのうち、当該文書が含むキーワード(「本塁打記録、本塁打、年間最多本塁打記録、号本塁打、46本、47号)を画面表示する。また、表示部15は、例えば、論文名が「A」である文書のテキストデータが表示されるセルに対応する行とクラスターBが表示されるセルに対応する列とが交差するセルの位置に、当該クラスターB中のキーワードのうち、当該文書が含むキーワード(「リーグ」、「米大リーグ」、「大リーグ」)を画面表示する。また、表示部15は、例えば、論文名が「A」である文書のテキストデータが表示されるセルに対応する行とクラスターCが表示されるセルに対応する列とが交差するセルの位置に、当該クラスターC中のキーワードのうち、当該文書が含むキーワード(「マグワイア」、「マグワイア一塁手」、「マーク」、「ロジャー」、「ソーサ」、「ソーサ外野

10

20

30

40

50

手」、「サミー」)を画面表示する。また、表示部15は、例えば、論文名が「A」である文書のテキストデータが表示されるセルに対応する行とクラスターDが表示されるセルに対応する列とが交差するセルの位置に、当該クラスターD中のキーワードのうち、当該文書が含むキーワード(「カージナルス」、「ヤンキース」)を画面表示する。また、表示部15は、例えば、論文名が「A」である文書のテキストデータが表示されるセルに対応する行とクラスターEが表示されるセルに対応する列とが交差するセルの位置に、当該クラスターE中のキーワードのうち、当該文書が含むキーワード(「試合」、「チーム」、「スタジアム」)を画面表示する。また、表示部15は、例えば、論文名が「A」である文書のテキストデータが表示されるセルに対応する行とクラスターFが表示されるセルに対応する列とが交差するセルの位置に、当該クラスターF中のキーワードのうち、当該文書が含むキーワード(「単独トップ」)を画面表示する。

10

【0118】

同様に、表示部15は、論文名が「B」という文書、論文名が「C」という文書についてのテキストデータに対応付けて、クラスター中のキーワードのうち各文書が含むキーワードを、クラスター毎に画面表示する。

【0119】

本発明の一実施形態によれば、例えば上記ステップS6の処理の後に、クラスターソート指定部16が、各クラスターのソート内容をデータソート部14に対して指示し、データソート部14が、指示されたソート内容に従って各クラスターをソートし、表示部15が、データソート部14によってソートされた文書データに対応付けて、データソート部14によってソートされたクラスター中のキーワードのうち当該ソートされた文書データを持つ文書が含むキーワードを、データソート部14によってソートされた各クラスター毎に画面表示するようにしてもよい。

20

【0120】

図5に示す表示画面を参照すれば、例えば、論文名が「A」という文書は、米大リーグの試合における本塁打記録に関する文書であることを把握することができる。従って、本発明によれば、文書の大凡の内容を一見して把握できるような表示を行うことが可能となる。

【0121】

なお、本発明は、コンピュータにより読み取られ実行されるプログラムとして実施することもできる。本発明を実現するプログラムは、コンピュータが読み取り可能な、可搬媒体メモリ、半導体メモリ、ハードディスクなどの適当な記録媒体に格納することができ、これらの記録媒体に記録して提供され、又は、通信インタフェースを介してネットワークを利用した送受信により提供されるものである。

30

【図面の簡単な説明】

【0122】

【図1】本発明のシステム構成の一例を示す図である。

【図2】データ表示処理フローの一例を示す図である。

【図3】各キーワードの頻度を示す図である。

【図4】各文書と各文書に含まれるキーワードとの対応情報を示す図である。

40

【図5】画面表示例を示す図である。

【符号の説明】

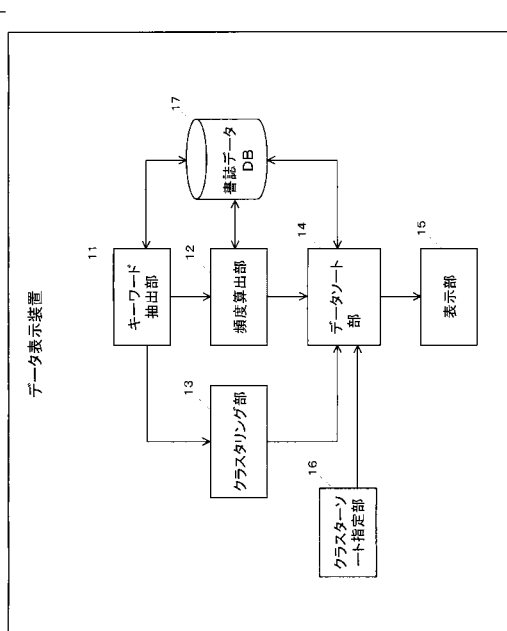
【0123】

- 1 データ表示装置
- 11 キーワード抽出部
- 12 頻度算出部
- 13 クラスタリング部
- 14 データソート部
- 15 表示部
- 16 クラスタソート指定部

50

17 書誌データDB

【図1】



【図2】

