

(19) 日本国特許庁(JP)

(12) 特許公報(B2)

(11) 特許番号

特許第3921540号
(P3921540)

(45) 発行日 平成19年5月30日(2007.5.30)

(24) 登録日 平成19年3月2日(2007.3.2)

(51) Int. Cl.

G06F 17/30 (2006.01)

F I

G06F 17/30 210D

G06F 17/30 170A

請求項の数 12 (全 22 頁)

<p>(21) 出願番号 特願2003-290929 (P2003-290929)</p> <p>(22) 出願日 平成15年8月8日(2003.8.8)</p> <p>(65) 公開番号 特開2005-63071 (P2005-63071A)</p> <p>(43) 公開日 平成17年3月10日(2005.3.10)</p> <p>審査請求日 平成15年8月8日(2003.8.8)</p> <p>特許法第30条第1項適用 平成15年2月13日鈴鹿サーキット フラワーガーデンホテルにおいて開催された「情報アクセスのためのテキスト処理」シンポジウムで発表</p>	<p>(73) 特許権者 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</p> <p>(74) 代理人 100130498 弁理士 佐野 禎哉</p> <p>(72) 発明者 野畑 周 東京都小金井市貫井北町4-2-1 独立行政法人通信総合研究所内</p> <p>(72) 発明者 井佐原 均 東京都小金井市貫井北町4-2-1 独立行政法人通信総合研究所内</p>
---	---

最終頁に続く

(54) 【発明の名称】 文書セット分類装置及びそのプログラム

(57) 【特許請求の範囲】

【請求項1】

任意に与えられた固有表現クラスの定義に基づき複数の文書の集合である文書セットに対して前記固有表現クラスの定義に基づいて得られる分類を付与するものであって、前記文書セットに含まれる文書の中に出現する固有表現を抽出するとともに、抽出した固有表現が属する固有表現クラスを、固有表現と固有表現が属する固有表現クラスとが列挙されたデータ群を参照して判定する固有表現抽出手段と、

前記文書セットの主題が単独の固有表現に関するものか複数の固有表現に関するものかを判断し、かつ、該固有表現が任意に与えられた前記固有表現クラスのうち何れの固有表現クラスに属するかを判断する判断手段と、

前記判断手段が下した判断に基づき、前記文書セットの主題に係る固有表現が単独であるか複数であるか、及び、該固有表現が属している前記固有表現クラスという2つの要素より規定される分類についての情報を出力する出力手段と

を具備し、

前記判断手段は、

判断対象である固有表現の頻度が予め定められた閾値以上となる場合はその単独の固有表現に関するものと判断し、当該固有表現がなく且つ判断対象である固有表現クラスの頻度がそれぞれに対して予め定められた閾値以上となる場合は複数の固有表現に関するものと判断するものであり、かつ、

前記固有表現が、前記固有表現抽出手段の判定結果に従って、その判定された固有表現ク

ラスに属すると判断するものである
ことを特徴とする文書セット分類装置。

【請求項2】

前記文書セットに含まれる文書の中に出現するクラスターを抽出するとともに、抽出したクラスターが関連する固有表現クラスを、クラスターとクラスターが関連する固有表現クラスとが列挙されたデータ群を参照して判定するクラスター抽出手段をさらに具備し、

前記判断手段は、

前記固有表現が、前記クラスター抽出手段の判定結果に従って、その判定された固有表現クラスに属すると判断するものである

請求項1記載の文書セット分類装置。

【請求項3】

任意に与えられた固有表現クラスの定義に基づき複数の文書の集合である文書セットに対して前記固有表現クラスの定義に基づいて得られる分類を付与するものであって、

前記文書セットに含まれる文書の中に出現するクラスターを抽出するとともに、抽出したクラスターが関連する固有表現クラスを、クラスターとクラスターが関連する固有表現クラスとが列挙されたデータ群を参照して判定するクラスター抽出手段と、

前記文書セットの主題が単独の固有表現に関するものか複数の固有表現に関するものかを判断し、かつ、該固有表現が任意に与えられた前記固有表現クラスのうち何れの固有表現クラスに属するかを判断する判断手段と、

前記判断手段が下した判断に基づき、前記文書セットの主題に係る固有表現が単独であるか複数であるか、及び、該固有表現が属している前記固有表現クラスという2つの要素より規定される分類についての情報を出力する出力手段と

を具備し、

前記判断手段は、

判断対象である固有表現の頻度が予め定められた閾値以上となる場合はその単独の固有表現に関するものと判断し、当該固有表現がなく且つ判断対象である固有表現クラスの頻度がそれぞれに対して予め定められた閾値以上となる場合は複数の固有表現に関するものと判断するものであり、かつ、

前記固有表現が、前記クラスター抽出手段の判定結果に従って、その判定された固有表現クラスに属すると判断するものである

ことを特徴とする文書セット分類装置。

【請求項4】

前記判断手段は、前記文書セットに含まれる複数の文書の各々が作成若しくは発表された時点に関する情報を参照し、これら複数の文書のうちの一定の割合以上のものが予め定められた期間内に作成若しくは発表されていることを条件として、前記記事セットの主題に係る固有表現が単独でありかつその属する固有表現クラスがイベント名クラスである旨の判断を下す請求項1、2または3記載の文書セット分類装置。

【請求項5】

与えられた文書の中に存在するキーワードを抽出し、一の文書のキーワードと他の文書のキーワードとの類似度を算出し、その類似度が閾値を超える場合にこれらの文書を同一の文書セットに割り当てることを通じて、複数の文書から少なくとも一の文書セットを生成し得る文書セット生成手段をさらに具備する請求項1、2、3または4記載の文書セット分類装置。

【請求項6】

請求項1、2、3、4または5記載の文書セット分類装置とともに用いられるものであって、

前記文書セット分類装置が出力する、前記文書セットに付与された分類についての情報を参照し、この分類に対応した要約アルゴリズムにより前記文書セットに含まれる複数の文書を単一の文書に要約する要約手段を具備する文書要約装置。

10

20

30

40

50

【請求項7】

請求項1、2、4または5記載の文書セット分類装置を構成するために用いられるものであって、コンピュータを、少なくとも、

複数の文書の集合である文書セットに含まれる文書の中に出現する固有表現を抽出するとともに、抽出した固有表現が属する固有表現クラスを固有表現と固有表現が属する固有表現クラスとが列挙されたデータ群を参照して判定する固有表現抽出手段、

文書セットの主題が単独の固有表現に関するものか複数の固有表現に関するものかを判断し、かつ、該固有表現が任意に与えられた前記固有表現クラスのうち何れの固有表現クラスに属するかを判断する判断手段、及び、

前記判断手段が下した判断に基づき、前記文書セットの主題に係る固有表現が単独であるか複数であるか、及び、該固有表現が属している前記固有表現クラスという2つの要素より規定される分類についての情報を出力する出力手段

として機能させ、

前記判断手段は、前記コンピュータを、

判断対象である固有表現の頻度が予め定められた閾値以上となる場合はその単独の固有表現に関するものと判断し、当該固有表現がなく且つ判断対象である固有表現クラスの頻度がそれぞれに対して予め定められた閾値以上となる場合は複数の固有表現に関するものと判断し、かつ、

前記固有表現が、前記固有表現抽出手段の判定結果に従って、その判定された固有表現クラスに属すると判断するように機能させる

ことを特徴とするプログラム。

【請求項8】

前記コンピュータを、さらに

前記文書セットに含まれる文書の中に出現するクラスターを抽出するとともに、抽出したクラスターが関連する固有表現クラスを、クラスターとクラスターが関連する固有表現クラスとが列挙されたデータ群を参照して判定するクラスター抽出手段として機能させ、

前記判断手段は、前記コンピュータを、

前記固有表現が、前記クラスター抽出手段の判定結果に従って、その判定された固有表現クラスに属するとの判断を実行するものである請求項7記載のプログラム。

【請求項9】

請求項3、4または5記載の文書セット分類装置を構成するために用いられるものであって、コンピュータを、少なくとも、

複数の文書の集合である文書セットに含まれる文書の中に出現するクラスターを抽出するとともに、抽出したクラスターが関連する固有表現クラスを、クラスターとクラスターが関連する固有表現クラスとが列挙されたデータ群を参照して判定するクラスター抽出手段と、

前記文書セットの主題が単独の固有表現に関するものか複数の固有表現に関するものかを判断し、かつ、該固有表現が任意に与えられた前記固有表現クラスのうち何れの固有表現クラスに属するかを判断する判断手段と、

前記判断手段が下した判断に基づき、前記文書セットの主題に係る固有表現が単独であるか複数であるか、及び、該固有表現が属している前記固有表現クラスという2つの要素より規定される分類についての情報を出力する出力手段と

として機能させ、

前記判断手段は、前記コンピュータを、

判断対象である固有表現の頻度が予め定められた閾値以上となる場合はその単独の固有表現に関するものと判断し、当該固有表現がなく且つ判断対象である固有表現クラスの頻度がそれぞれに対して予め定められた閾値以上となる場合は複数の固有表現に関するものと判断し、かつ、

前記固有表現が、前記クラスター抽出手段の判定結果に従って、その判定された固有表

10

20

30

40

50

現クラスに属すると判断するように機能させる
ことを特徴とする文書セット分類装置。

【請求項 10】

前記判断手段は、前記文書セットに含まれる複数の文書の各々が作成若しくは発表された時点に関する情報を参照し、これら複数の文書のうちの一定の割合以上のものが予め定められた期間内に作成若しくは発表されていることを条件として、前記記事セットの主題に係る固有表現が単独でありかつその属する固有表現クラスがイベント名クラスである旨の判断を下す請求項 7、8 または 9 記載のプログラム。

【請求項 11】

さらにコンピュータを、与えられた文書の中に存在するキーワードを抽出し、一の文書のキーワードと他の文書のキーワードとの類似度を算出し、その類似度が閾値を超える場合にこれらの文書を同一の文書セットに割り当てることを通じて、複数の文書から少なくとも一の文書セットを生成し得る文書セット生成手段としても機能させる請求項 7、8、9 または 10 記載のプログラム。

【請求項 12】

請求項 6 記載の文書要約装置を構成するために用いられるものであって、コンピュータを、少なくとも、文書セット分類装置が出力する、文書セットに付与された分類についての情報を参照し、この分類に対応した要約アルゴリズムにより前記文書セットに含まれる複数の文書を単一の文書に要約する要約手段として機能させるプログラム。

【発明の詳細な説明】

【技術分野】

【0001】

本発明は、複数の文書を包含する文書セットに分類を付与するための分類装置及びそのプログラムに関する。

【背景技術】

【0002】

複数文書の自動要約は、要約の研究において近年関心の高まっている分野である。米国の Document Understanding Conference (DUC) や日本の Text Summarization Challenge (TSC) でも、要約システムの評価を行う課題として複数文書の要約が対象に加えられている。複数文書要約とは、単一の主題について収集された複数の文書を含む文書セットを単一の文書に要約することである。より具体的に述べると、ある事件の始まりから終わりまでの一連の報告や、特定個人の行動・発言の内容、各地で発生した地震の被害状況等の主題に沿って収集された複数の文書より、当該主題に関する要約を生成することである。

【0003】

要約の精度を向上させるためには、文書セットがもつ主題を正しく把握し、それに応じて適切な要約手法、出力形式を選択する必要があると考えられる。複数文書要約の観点から文書セットを分類する先行研究として、コロンビア大学の McKeown 等によるものがある (非特許文献 1 を参照)。McKeown 等は、複数の新聞記事を包含する記事セットに付与すべき分類として、

(A) Single-Event (特定の地域・期間に限定された単独の出来事に関する記事セット)

(B) Person-centered (特定人物にまつわる出来事を記述した記事セット)

(C) Multi-Event (異なる地域・期間にまたがった複数の出来事に関する記事セット。

大抵は行動主体も異なる)

(D) Other (上記の 3 分類に当てはまらない、互いに漠然と関連している記事セット)

の 4 種類を定義した。そして、記事セットを分類する際の手がかりとして、記事セット中の全記事間のタイムスパン、同日に掲載された記事の割合、大文字で始まる語の頻度、he、she 等の人称代名詞の頻度、を用いている。

【非特許文献 1】 K. R. McKeown and R. Barzilay and D. Evans and V. Hatzivassilou and M. Yen Kan and B. Schiffman and S. Teufel, [online], "Columbia Multi-Docu

10

20

30

40

50

ment Summarization: Approach and Evaluation”, Online Proceedings of DUC2001 <http://www-nlpir.nist.gov/projects/duc/pubs/2001papers/columbia#redo.pdf>

【発明の開示】

【発明が解決しようとする課題】

【0004】

McKeown等の分類は、要約の対象となる記事セットによく見られる性質を効率よく分類している。しかしながら、幾つかの問題も抱えている。即ち、

・Otherに分類される記事セットが多くなる。これらに対しては他に適切な分類があるのではないかと考えられる。

・分類を判定するアルゴリズムに用いられている手がかりのうち、大文字で始まる語の頻度及びhe、she等の人称代名詞の頻度は英語に特化したものである。より一般的に記事セットの分類を行うためには、用いるべき手がかりを考慮する必要がある。

・上記非特許文献1で分類対象とされた記事データは、複数記事要約の評価ワークショップで使用するために作成されたものである。そのために、一般的な記事セットと比較して整えられているか、あるいは偏りが生じている。

・McKeown等の要約システムでは、実際に複数記事要約を行うときにMulti-EventとOtherとを同一視して同じ要約手法を適用しており、Multi-EventとOtherとを区別した意義が失われている。

【0005】

以上に鑑みてなされた本発明は、特定言語の特性に依存せずなおかつ分類の網羅性を高めたより好適な分類を定義した上で、これに基づいた分類を文書セットに付与できる分類装置を提供するものである。

【課題を解決するための手段】

【0006】

上述した課題を解決すべく、本発明では、任意に与えられた固有表現クラスの定義に基づき複数の文書の集合である文書セットに対し前記固有表現クラスの定義に基づいて得られる分類を付与するものとして、図1に示すように、前記文書セットに含まれる文書の中に出現する固有表現を抽出するとともに、抽出した固有表現が属する固有表現クラスを、固有表現と固有表現が属する固有表現クラスとが列挙されたデータ群を参照して判定する固有表現抽出手段（図示省略）と、前記文書セットの主題が単独の固有表現（Named Entity）に関するものか複数の固有表現に関するものかを判断し、かつ、該固有表現が何れの固有表現クラスに属するかを判断する判断手段101と、前記判断手段101が下した判断に基づき、前記文書セットの主題に係る固有表現が単独であるか複数であるか、及び、該固有表現が属している前記固有表現クラスという2つの要素より規定される分類についての情報を出力する出力手段102とを具備する文書セット分類装置を構成した。

【0007】

本発明では、文書セットがもつ主題に関連する固有表現が属する固有表現クラスに基づいて分類を定義し、本発明に係る分類装置がこの定義された分類を文書セットに付与するものとした。このようなものであれば、特定言語に依存する度合いが低減するとともに、分類の網羅性が高まる、言い換えるならば文書セットがOtherに分類される可能性が小さくなる。

【0008】

ここで、固有表現とは、文書中に含まれる人名、組織名等の固有名詞や日付、金額等の数値表現その他の情報抽出の要素となる表現を言う。固有表現は、情報として重要でありかつその内容を示す表現がほぼ一意に定まるものである。固有表現が属する固有表現クラスの定義は種々考えられるが、一態様として、人名（Person）、組織名（Organization）、地名（Location）、施設名（Facility）、固有物名（Product（製品名、法律名等））、イベント名（Event）の6種を内包するクラス定義を採用することができる。この固有表現クラス定義を採用した場合、文書セットに付与すべき分類は以下の通りとなる。即ち、

10

20

30

40

50

- (A) Single-Person (単一人物に関する文書セット)
- (B) Single-Organization (単一組織に関する文書セット)
- (C) Single-Location (単一地域に関する文書セット)
- (D) Single-Facility (単一施設に関する文書セット)
- (E) Single-Product (単一固有物に関する文書セット)
- (F) Single-Event (単一イベントに関する文書セット)
- (G) Multi-Person (複数人物に関する文書セット)
- (H) Multi-Organization (複数組織に関する文書セット)
- (I) Multi-Location (複数地域に関する文書セット)
- (J) Multi-Facility (複数施設に関する文書セット)
- (K) Multi-Product (複数固有物に関する文書セット)
- (L) Multi-Event (複数イベントに関する文書セット)
- (M) Other (その他)

10

但し、固有表現クラスの定義ひいては文書セットの分類の定義はこれに限定されない。よって、例えば、動物名クラス (Single/Multi-Animal) 等を追加することができ、人名、動物名等を一のクラス、いわば行動主体を表現するクラスである「主体名クラス」にまとめることもできる。

【 0 0 0 9 】

文書セット分類装置における前記判断手段 1 0 1 は、ここで前記判断手段は、判断対象である固有表現の頻度が予め定められた閾値以上となる場合はその単独の固有表現に関するものと判断し、当該固有表現がなく且つ判断対象である固有表現クラスの頻度がそれぞれに対して予め定められた閾値以上となる場合は複数の固有表現に関するものと判断を実行し、かつ、前記固有表現が、前記固有表現抽出手段の判定結果に従って、その判定された固有表現クラスに属するとの判断を実行する。但し、ここに言う頻度は、文書セットに含まれる複数の文書中に出現する特定の固有表現等の出現頻度には限られず、特定の固有表現等が出現する文書の数 (文書セットにおける文書の頻度) であることがある。

20

【 0 0 1 0 】

加えて、文書セット分類装置は、文書セットに含まれる文書の中に出現するクラスタームを抽出するとともに、抽出したクラスタームが関連する固有表現クラスを、クラスタームとクラスタームが関連する固有表現クラスとが列挙されたデータ群を参照して判定するクラスターム抽出手段をさらに具備する構成として、前記判断手段 1 0 1 が、前記固有表現が、前記クラスターム抽出手段の判定結果に従って、その判定された固有表現クラスに属するとの判断を実行するものであってもよい。或いは文書セット分類装置を、前述した固有表現抽出手段、判断手段、出力手段を備える構成に代えて、クラスターム抽出手段、判断手段、出力手段を備える構成とすることもできる。なお、クラスタームとは、特定の固有表現クラスに関連の強い名詞または複合名詞のことである。例えば、「首相」等の役職名は人名クラスのクラスタームであり、「地震」等の名詞はイベント名クラスのクラスタームである。因みに、クラスタームは、固有表現そのものとは異なり、一般名詞である。

30

【 0 0 1 1 】

また、特に、文書セットに含まれる各文書の作成時または発表時が判明しているような場合において、前記判断手段 1 0 1 が、前記文書セットに含まれる複数の文書の各々の作成若しくは発表された時点に関する情報を参照し、これら複数の文書のうちの一定の割合以上のものが予め定められた期間内に作成若しくは発表されていることを条件として、前記記事セットの主題に係る固有表現が単独でありかつその属する固有表現クラスがイベント名クラスである旨の判断を下すものとしてもよい。このとき、当該文書セットには Single-Event の分類が付与される。

40

【 0 0 1 2 】

図 2 に示すように、上記の文章セット分類装置が、与えられた文書の中に存在するキーワードを抽出し、一の文書のキーワードと他の文書のキーワードとの類似度を算出し、そ

50

の類似度が閾値を超える場合にこれらの文書を同一の文書セットに割り当てることを通じて、複数の文書から少なくとも一の文書セットを生成し得る文書セット生成手段103をさらに具備するものであれば、与えられた複数の文書を一または複数の文書セットに仕分けしこれに分類を付与するまでの処理を一括に実行可能となる。このものは、与えられた文書を基に一または複数の要約を自動生成するシステムを構築するために有用となる。

【0013】

さらに、図3に示すように、上記の文書セット分類装置が文書セットに対して付与した分類を参照し、この分類に対応した要約アルゴリズムを選択して前記文書セットに含まれる複数の文書を要約する要約手段201を具備する文書要約装置を構成することで、より適切に複数文書の要約を実行することが可能となる。

10

【発明の効果】

【0014】

以上に詳述した本発明によれば、特定言語の特性に依存せず、かつ分類の網羅性を高めたより好適な分類を文書セットに付与し得る。

【発明を実施するための最良の形態】

【0015】

以下、本発明の一実施形態を、図面を参照して説明する。はじめに、本発明における分類の定義及びその妥当性について述べる。ここでは、分析対象とする文書セットとして複数の新聞記事を包含する記事セットを実験的に生成し、これを分析して固有表現クラスを基にした分類を定義する。そして、この分類を、テストデータとなる別の日本語新聞記事セット、及び、DUC2001で使用された英語新聞記事セットに適用することにより、分類の定義の妥当性を検証する。

20

【0016】

記事セットの偏りを避けるため、日本語新聞記事コーパスから無作為に一の記事を抽出し、その記事に類似する記事を情報検索システムを使用して収集して、記事セットを生成した。具体的な手順は以下の通りである。

- (1) 日本語新聞記事コーパスに含まれる記事から無作為に記事の一つを選択する
- (2) 選択した記事からキーワード列を抽出する；キーワード列は、既知の日本語形態素解析ソフトウェア（黒橋禎夫、長尾真 日本語形態素解析システム Juman version 3.61. 京都大学, 1999.）を用いた形態素解析結果より時相名詞、副詞的名詞を除いた名詞のうち頻度2以上のものとした
- (3) 抽出したキーワード列を用いて、記事間の類似度を求める；各記事について上記(2)と同様にキーワード列を抽出し、キーワード同士の類似度をDice's coefficient (Diceの係数)を用いて求めた
- (4) 類似した記事を取り出す；同一の記事以外でDiceの係数が所定値（例えば、0.5）以上となる記事を類似記事と見なして取り出した。

30

【0017】

以上の記事セット生成を複数回（例えば、50回）繰り返して得られた複数（50）記事セットのうち、3以上の記事を含む記事セットを選出すると26記事セットとなった。これらから、さらに記事セットの内容がほぼ同じと考えられるものを省くと、19記事セットが残った。これらの19記事セットの主題を、図4に示す。

40

【0018】

図4の分析結果に示されるように、Singleとは、記事セット中のほとんどの記事が単一のイベントや人名、組織名等の固有表現について記述してあるものである。他方、Multiとは、記事セット中の記事が複数の相異なるイベントや人名、組織名等の固有表現について記述してあるものである。

【0019】

本発明では、文書セットたる記事セットの主題を、固有表現クラスのうち一つを選択することを通じて分類した。固有表現クラスの定義としては、拡張されたクラス定義（S. Sekine, K. Sudo and C. Nobata, "Extended named entity hierarchy", In Proceeding

50

s of the LREC-2002 conference, 2002.) を採用した。この定義は階層的であるが、最上位の階層を用いてほぼ全ての記事セットの主題を分類することができた。固有表現クラスを割り当てず、Otherに分類した記事セットは一つだけである(図4における記事セット No. 14)。この記事セットに固有表現クラスを割り当てるとするならば動物名(Animal)クラスとなるが、動物名クラスに分類される記事セットはそれほど多くはないと考えられるので、Otherとした。上述したように固有表現クラスを選択的に割り当てた結果、既に言及した13種類の分類が定義された。因みに、図4に示した分析結果には地名(Location)クラスの分類に対応する記事セットが存在しないが、特定の国や地域に関する記事セットは存在し得るし重要でもあると考えられるので、記事セットの分類の定義に含めた。

10

【0020】

続いて、定義した分類に基づいてテストデータの分類を行う。ここで使用するテストデータは、日本語新聞記事コーパスから先に述べた方法と同様にして作成される。そして、テストデータとして作成された20の記事セットについて、二人の被験者が独立に分類を付与した。原則として、被験者は各記事セットに一つの分類を割り当てたが、幾つかの記事セットについては複数の分類が可能であると判断して二つの分類を割り当てた。二人の被験者によって割り当てられた各分類の数を、図5に示す。被験者間の一致率は、被験者が最初に選択した分類同士を比較した場合には55%、二番目の分類までを含めた場合には85%であった。Otherに分類された記事セットはなかった。被験者が二つの分類を割り当てた記事セットの数は、それぞれ6と5であった。

20

【0021】

被験者による分類結果を基に、テストデータの分類の正解データを作成した。被験者間で共通の分類となった17記事セットではその割り当てられた分類を正解とする一方、被験者間で分類が分かれた3記事セットについては被験者同士の討論により正解の分類を決定した。なお、この正解データを用いて、後述する文書セット分類装置による自動分類の実験結果の評価を行う。

【0022】

この分類の定義が他言語においても妥当であることを示すために、DUC2001の複数記事要約タスクで用いられたトレーニングデータである英語新聞記事セットに対しても分類の付与を試みた。先と同様に、二人の被験者が独立に記事セットに対して最適と判断した分類を付与した。被験者は各記事セットに対して一つないし二つの分類を付与した。二人の被験者によって割り当てられた各分類の数を、図6に示す。被験者が二つの分類を割り当てた記事セットはそれぞれ4セットあった。Otherの分類は、被験者の一人が二番目の分類として付与した1セットのみであった。被験者間の一致率は、被験者が最初に選択した分類同士を比較した場合に80%、二番目の分類までを含めた場合は93.3%であった。この一致率は、日本語記事のテストデータに対するものよりも高い。これは、DUC2001のデータが日本語記事テストデータよりも整えられており、記事セットの主題が意図的に選択されているからであると考えられる。

30

【0023】

以降、定義された分類を記事セットに付与するための文書セット分類装置について詳述する。本実施形態における文書セット分類装置は、コンピュータ1に所定のプログラムをインストールすることで構成されるものである。コンピュータ1は、例えば、図7に示すように、プロセッサ1a、メインメモリ1b、ハードディスクドライブに代表される補助記憶デバイス1c等のハードウェア資源が、コントローラ1d(即ち、いわゆるシステムコントローラ、I/Oコントローラ等)により制御され連携して動作するものである。また、図示しないが、電気通信回線を介して外部とのデータ授受を行うための通信デバイス、ユーザによる操作入力を受け付けるキーボードやポインティングデバイス等の入力デバイス、情報を画像ないし映像として表示するディスプレイ及びこのディスプレイに映像信号を送出するため表示制御デバイス(いわゆるグラフィクスチップ等)等を実装することを妨げない。

40

50

【 0 0 2 4 】

通常、プロセッサ 1 a によって実行されるべきプログラムが補助記憶デバイス 1 c に格納されており、プログラムの実行の際には補助記憶デバイス 1 c からメインメモリ 1 b に読み込まれ、プロセッサ 1 a によって解読される。そして、該プログラムに従い上記のハードウェア資源を作動して、少なくとも、判断手段 1 0 1、出力手段 1 0 2 としての機能を発揮するようにしている。

【 0 0 2 5 】

判断手段 1 0 1 は、複数の記事（文書）を包含してなる記事セット（文書セット）の主題が単独の固有表現に関するものか複数の固有表現に関するものかを判断し、かつ、該固有表現が何れの固有表現クラスに属するかを判断する。入力として与えられる記事セットの要素である記事のデータは、通常、メインメモリ 1 b または補助記憶デバイス 1 c の所要の記憶領域に予め格納されている。よって、プロセッサ 1 a が、プログラムに従い、メインメモリ 1 b または補助記憶デバイス 1 c に格納されている記事のデータを読み込み、これを基に記事セットに付与すべき分類についての判定を行う。

10

【 0 0 2 6 】

出力手段 1 0 2 は、前記判断手段 1 0 1 が下した判断より、前記記事セットに付与すべき分類についての情報、即ち、当該記事セットの主題に係る固有表現が単独であるか複数であるか及び該固有表現が属している固有表現クラスという 2 つの要素より規定される分類についての情報を出力する。情報の出力の態様としては、ディスプレイの画面への表示、プリンタ（図示せず）を使用したプリントアウト、通信デバイス及び電気通信回線を介した外部のコンピュータへの送信、メインメモリ 1 b または補助記憶デバイス 1 c への書き込み、その他を挙げることができる。出力手段 1 0 2 の具体的構成は、記事セットに付与すべき情報の出力態様に応じたものとなる。

20

【 0 0 2 7 】

文書セット分類装置による分類の自動付与では、記事セットに含まれる複数の記事中に出現する単語や固有表現クラスの出現頻度と、記事頻度（記事セットにおいて所要の単語、固有表現クラス等が出現した記事の数。情報検索等で用いられる i d f 値とは異なる）とを手がかりとして利用する。よって、コンピュータ 1 が、図 8 に示す固有表現抽出手段としての機能をも発揮し得ることが望ましい。固有表現抽出手段は、記事セットに含まれる記事の中に出現する固有表現を抽出するとともに、抽出した固有表現が属する固有表現クラスを判定する。固有表現抽出手段は、例えば、記事中の文章を単語に切り分けて品詞の付与を行う形態素解析ソフトウェア 1 0 4 a と、形態素解析ソフトウェア 1 0 4 a による解析結果を参照して記事中に出現する固有表現を列挙しこれを固有表現クラス分けする固有表現抽出ソフトウェア 1 0 4 b とを用いて構成できる。形態素解析、固有表現抽出の一例を、図 8 に示している。図示例では、パターンベースのシステムで固有表現抽出を行っている。即ち、プロセッサ 1 a が、プログラムに従い、入力として与えられた記事を形態素解析し、得られた形態素解析済み記事を固有表現リスト 1 0 4 c（固有表現及びその属する固有表現クラスが列挙されたデータ群。通常、メインメモリ 1 b または補助記憶デバイス 1 c の所要の記憶領域に格納されている）に照らし合わせることで、記事中の固有表現を全て抽出する。しかる後、複数の固有表現が入れ子関係となっているもの（例えば、組織名クラスに属する固有表現「吉本工業」の中に、さらに人名クラスに属する固有表現「吉本」が存在）が存在しているときにはより文字列の長い固有表現を優先的に認定（即ち、「吉本」ではなく「吉本工業」という固有表現と認定）して固有表現を一意に決定し、その結果を出力する。判断手段 1 0 1 は、この固有表現抽出手段による出力を参照して、記事セットの分類を行うことができる。

30

40

【 0 0 2 8 】

また、固有表現に加え、記事中に出現するクラスターも分類の手がかりとして用いることができる。クラスターとは、特定の固有表現クラスに関連の強い名詞または複合名詞のことである。例示すると、「首相」等の役職名は人名クラスのクラスターであり、「地震」等の名詞はイベント名クラスのクラスターである。よって、コンピュータ 1 が

50

、図9に示すクラスターム抽出手段としての機能をも発揮し得ることが望ましい。クラスターム抽出手段は、記事セットに含まれる記事の中に出現するクラスタームを抽出するとともに、抽出したクラスタームが関連する固有表現クラスを判定する。クラスターム抽出手段は、例えば、記事中の文章を単語に切り分けて品詞の付与を行う形態素解析ソフトウェア104aと、形態素解析ソフトウェア104aによる解析結果を参照して記事中に出現するクラスタームを列挙しこれを固有表現クラス分けするクラスターム抽出ソフトウェア105aとを用いて構成できる。形態素解析、クラスターム抽出の一例を、図9に示している。図示例もまた、上記の固有表現抽出の例と同様に、パターンベースのシステムでクラスターム抽出を行っている。即ち、プロセッサ1aが、プログラムに従い、入力として与えられた記事を形態素解析し、得られた形態素解析済み記事をクラスタームリスト105b（クラスターム及びその関連する固有表現クラスが列挙されたデータ群。通常、メインメモリ1bまたは補助記憶デバイス1cの所要の記憶領域に格納されている）に照らし合わせることで記事中のクラスタームを全て抽出し、その結果を出力する。判断手段101は、このクラスターム抽出手段による出力を参照して、記事セットの分類を行うことができる。因みに、クラスタームのリスト105bは、既存のシソーラスと人手による収集結果とから作成することができる。発明者が実験的に作成し使用しているクラスタームの数は約16000語である。

【0029】

固有表現、クラスタームの他、記事の掲載日付や記事間のタイムスパン等もまた、記事セットの分類のための手がかりとして用いることが可能である。例えば、ある記事セットに含まれるほとんどの記事が同日かまたは所定の短い期間以内に掲載されたものであるならば、その記事セットをSingle-Eventに分類できる可能性が高い。

【0030】

本実施形態における判断手段101が実行する判断のアルゴリズムに関して、図10ないし図14のフローチャートを参照して詳述する。本実施形態において、判断手段101は、下記の4つのアルゴリズムにより記事セットの分類を行う。まず、判断手段101は、第一のアルゴリズムに従い、入力として与えられた記事セットがSingle-Eventに分類されるか否かを判断する。この判断は、記事セットに含まれる複数の記事の各々が作成若しくは発表された時点に関する情報を参照して下される。記事セットに含まれる記事のうち一定の割合以上のもの（「一定の割合」の設定値により、記事セットに含まれる記事のうち一部または全部の何れかが該当する）が、所定期間内に作成、公開、発表、掲載等されたものであるならば、判断手段101は当該記事セットをSingle-Eventに分類する。より具体的には、記事セットに含まれる記事が新聞記事である場合に、その大半が同日あるいは所定の短い期間内に掲載されたものであるならば当該記事セットにSingle-Eventの分類を付与する旨の判断を下す。このときの処理の手順を、図10に示している。判断手段101は、入力として与えられた記事セットに含まれている各記事の掲載日付に関する情報（この情報は、例えば、記事データに関連づけてメインメモリ1bまたは補助記憶デバイス1cの所要の記憶領域に格納されている）を参照して（ステップS101）最も記事頻度の高い掲載日を確認し、この日に掲載された記事の数を計数する（ステップS102）。そして、この日に掲載された記事数またはこの日に掲載された記事数の記事セットに含まれる全記事数に対する割合が、予め設定された閾値 T_0 を上回っているならば（ステップS103）、与えられた記事セットをSingle-Eventに分類する（ステップS106）。また、記事セットに含まれる全記事の掲載日のタイムスパンを確認し（ステップS104）、このタイムスパンが予め設定された閾値 T_0 を下回っているならば（ステップS105）、与えられた記事セットをSingle-Eventに分類する（ステップS105）。つまり、特定の日に掲載された記事の割合、及び、記事間のタイムスパンの最大値の二つのパラメータを材料として判断を行う。記事セットをSingle-Eventに分類できなかったときには、第二のアルゴリズムに移行する。なお、記事中に出現する日付表現を参照することで第一のアルゴリズムを実行することを妨げない。

【0031】

次に、判断手段101は、第二のアルゴリズムに従い、与えられた記事セットがSingle-class（何れかの固有表現クラス）に分類されるかどうか判断する。この判断は、記事頻度の高い固有表現について、その出現頻度、記事頻度を計数することで下される。特定の固有表現が記事セットに含まれる多くの記事にわたって頻繁に出現するならば、判断手段101はその固有表現を記事セットの主題を表すものと見なし、その固有表現が属する固有表現クラスを記事セットの分類の要素classとする。このときの処理の手順を、図11に示している。判断手段101は、与えられた記事セットに対して前記固有表現抽出手段が行った固有表現抽出処理の結果出力を参照し（ステップS201）、記事セットに含まれる記事中に出現する固有表現及び固有表現が属する固有表現クラスのそれぞれについて、出現頻度を計数する（ステップS202）。かつ、記事中に出現する各固有表現についてその記事頻度を計数して、当該記事セットにおいて最も記事頻度の高い固有表現を選出する（ステップS203）。ここで、記事頻度が等しい複数の固有表現が存在する場合には、例えば出現頻度のより高い固有表現を選択する。しかして、選出された固有表現の記事頻度（または、選出された固有表現の記事頻度の記事セットに含まれる全記事数に対する割合）が予め設定された閾値 T_0 を上回っているならば（ステップS204）、この固有表現が属する固有表現クラスを記事セットの分類の要素classとすることができる。なお、記事セットをSingle-classに分類するに際し、さらなる判断処理を付加することを妨げない。即ち、図11に示しているように、選出された固有表現の出現頻度/選出された固有表現が属する固有表現クラスの出現頻度、の比が予め定められた閾値 T_w を上回っていることを条件として（ステップS207）、選出された固有表現が属する固有表現クラスを記事セットの分類の要素classとする（ステップS208）ものとしてもよい。但し、本実施形態では、トレーニングデータを調査した実験結果から、イベント名クラスに関しては他の固有表現クラスに優先して判断するものとした。従って、選出された固有表現がイベント名クラスに属するときには（ステップS205）、与えられた記事セットをそのままSingle-Eventに分類する（ステップS206）ようにしている。記事セットをSingle-classに分類できなかつたときには、第三のアルゴリズムに移行する。

【0032】

続いて、判断手段101は、第三のアルゴリズムに従い、与えられた記事セットがMulti-class（何れかの固有表現クラス）に分類されるかどうか判断する。この判断は、記事頻度の高い固有表現クラスについて、その出現頻度、記事頻度を計数することで下される。特定の固有表現クラスに属する固有表現が記事セットに含まれる多くの記事にわたって頻繁に出現するならば、判断手段101はその固有表現クラスを記事セットの主題を表すものと見なし、記事セットの分類の要素classとする。このときの処理の手順を、図12に示している。判断手段101は、与えられた記事セットに対して前記固有表現抽出手段が行った固有表現抽出処理の結果出力を参照し（ステップS201）、記事セットに含まれる記事中に出現する固有表現及び固有表現が属する固有表現クラスのそれぞれについて、出現頻度を計数する（ステップS202。これらの処理は、既に第二のアルゴリズムにおいて実行されている）。かつ、記事中に出現する固有表現が属する各固有表現クラスについてその記事頻度を計数して、当該記事セットにおいて最も記事頻度の高い固有表現クラスを選出する（ステップS301）。ここで、記事頻度が等しい複数の固有表現クラスが存在する場合には、例えば出現頻度のより高い固有表現クラスを選択する。しかして、選出された固有表現クラスの記事頻度（または、選出された固有表現クラスの記事頻度の記事セットに含まれる全記事数に対する割合）が予め設定された閾値 T_0 を上回っているならば（ステップS302）、この固有表現クラスを記事セットの分類の要素classとすることができる。なお、記事セットをMulti-classに分類するに際し、さらなる判断処理を付加することを妨げない。即ち、図12に示しているように、選出された固有表現クラスの出現頻度/全固有表現クラス（全固有表現）の出現頻度、の比が予め定められた閾値 T_0 を上回っていることを条件として（ステップS305）、選出された固有表現クラスを記事セットの分類の要素classとする（ステップS306）ものとしてもよい。但し、本実施形態では、イベント名クラスに関しては他の固有表現クラスに優先して判断するも

10

20

30

40

50

のとし、選出された固有表現クラスがイベント名クラスに属するときには（ステップS 303）、与えられた記事セットをそのままMulti-Eventに分類する（ステップS 304）ようにしている。加えて、固有表現毎に相異なる閾値 T_c を設定することを妨げない。例えば、選出された固有表現クラスの出現頻度/全固有表現クラスの出現頻度の値と比較される閾値 T_c について、選出された固有表現クラスが地名クラス、組織名クラス、人名クラスの何れかである場合にはより厳しい即ちより大きい閾値 T_{c1} を適用し、選出された固有表現クラスが上記以外である場合にはより緩い即ちより小さい閾値 T_{c2} （ $T_{c1} > T_{c2}$ ）を適用することができる。記事セットをMulti-classに分類できなかったときには、第四のアルゴリズムに移行する。

【0033】

第三のアルゴリズムまでの過程で記事セットに付与すべき適切な分類を見出せなかった場合、判断手段101は、第四のアルゴリズムに従い、付与すべき分類を検討する。第四のアルゴリズムは、第二のアルゴリズムないし第三のアルゴリズムを、固有表現でなくクラスターを対象として実行するものと言える。即ち、記事中に出現するクラスター間の頻度またはクラスターが関連する固有表現クラスの頻度のうち少なくとも一方を材料として、与えられた記事セットの主題に係る固有表現が属する固有表現クラスの判断を下す。なお、特定のクラスターが記事セット中の多くの記事にわたって頻繁に出現していても、当該記事セットに割り当てべき分類はSingle-classでなくMulti-classとすることが望ましい。これは、クラスターは固有名詞ではなく一般名詞であって、複数種の固有表現を指示し得るものであることによる。このときの処理の手順を、図13及び図14に示している。判断手段101は、与えられた記事セットに対して前記クラスター抽出手段が行ったクラスター抽出処理の結果出力を参照し（ステップS 401）、記事セットに含まれる記事中に出現するクラスター及びクラスターが関連する固有表現クラスのそれぞれについて、出現頻度を計数する（ステップS 402）。かつ、記事中に出現する各クラスターについてその記事頻度を計数して、当該記事セットにおいて最も記事頻度の高いクラスターを選出する（ステップS 403）。ここで、記事頻度が等しい複数のクラスターが存在する場合には、例えば出現頻度のより高いクラスターを選択する。しかして、選出されたクラスターの記事頻度（または、選出されたクラスターの記事頻度の記事セットに含まれる全記事数に対する割合）が予め設定された閾値 T_d を上回っているならば（ステップS 404）、このクラスターが関連する固有表現クラスを記事セットの分類の要素classとすることができる。なお、記事セットをMulti-classに分類するに際し、さらなる判断処理を付加することを妨げない。即ち、図13に示しているように、選出されたクラスター間の出現頻度/選出されたクラスターが関連する固有表現クラスの出現頻度、の比が予め定められた閾値 T_w を上回っていることを条件として（ステップS 407）、選出された固有表現が属する固有表現クラスを記事セットの分類の要素classとする（ステップS 408）ものとしてもよい。但し、イベント名クラスに関しては他の固有表現クラスに優先して判断するものとし、選出されたクラスターがイベント名クラスに関連するものであるときには（ステップS 405）、与えられた記事セットをそのままMulti-Eventに分類する（ステップS 406）ようにしている。上記に加えて、記事中に出現するクラスターが関連している各固有表現クラスについてその記事頻度を計数し、当該記事セットにおいて最も記事頻度の高い固有表現クラスを選出する（ステップS 409）。記事頻度が等しい複数の固有表現クラスが存在する場合には、例えば出現頻度のより高い固有表現クラスを選択する。しかして、選出された固有表現クラスの記事頻度（または、選出された固有表現クラスの記事頻度の記事セットに含まれる全記事数に対する割合）が予め設定された閾値 T_d を上回っているならば（ステップS 410）、この固有表現クラスを記事セットの分類の要素classとすることができる。記事セットをMulti-classに分類するに際しては、さらなる判断処理を付加することができる。即ち、図14に示しているように、選出された固有表現クラスの出現頻度/全固有表現クラス（全クラスター）の出現頻度、の比が予め定められた閾値 T_c を上回っていることを条件として（ステップS 413）、選出された固有表現クラスを記事セットの分類の要素classと

10

20

30

40

50

する（ステップS 4 1 4）ものとする。但し、イベント名クラスに関しては他の固有表現クラスに優先して判断するものとし、選出された固有表現クラスがイベント名クラスに属するときには（ステップS 4 1 1）、与えられた記事セットをそのままMulti-Eventに分類する（ステップS 4 1 2）ようにしている。なお、ここでも、第三のアルゴリズムと同様、固有表現毎に相異なる閾値 T_{c1} 、 T_{c2} を設定することが許容される。

【0034】

上記の全てのアルゴリズムを用いても分類を付与できなかった場合、判断手段101は、予め定められたデフォルトの分類を当該記事セットに付与する（ステップS 4 1 5）。デフォルトの分類は、例えば、Multi-EventまたはOtherとする。

【0035】

上述のテストデータに対し、本実施形態の文書セット分類装置を使用して分類を付与する自動分類実験を行った結果について述べる。なお、アルゴリズム中の各閾値の決定は、ここではトレーニングデータを基に人手で行う。各閾値の設定は、 $T_a = 0.33$ 、 $T_s = 150$ 、 $T_d = 0.90$ 、 $T_w = 0.40$ 、 $T_{c1} = 0.80$ 、 $T_{c2} = 0.40$ とした。但し、閾値の大きさがここに示す値に限られないことは言うまでもない。テストデータに対する自動分類実験の結果の評価を、図15に示す。図15には、被験者による分類付与の結果の評価及びベースラインをも示した。被験者の評価は、各被験者が付与した分類と正解との比較評価である。両被験者の正解に対する評価は両被験者間の一致率55%よりも高いが、これは分類の正解が両被験者による分類付与結果を総合して作成されたためである。ベースラインは、トレーニングデータにおいて最も頻度の高い分類（この実験では、Single-Event）の記事セットのテストデータにおける数（及び、占める割合）である。文書セット分類装置について、「一致」の値は文書セット分類装置が出力した分類が正解に一致した記事セット数（及び、割合）を示し、「部分一致」の値は文書セット分類装置が出力した分類が被験者によって付与された分類の何れかに一致した数（及び、割合）を示す。被験者が複数の分類を付与した記事セットに関してはその双方を含む。被験者について、「一致」の値は被験者が最初に与えた分類が正解の分類に一致した記事セット数（及び、割合）を示し、「部分一致」の値は被験者が二番目に与えた分類も含めて正解の分類に一致した記事セット数（及び、割合）を示す。

【0036】

文書セット分類装置は、20記事セットのうち9つを正しく分類し、さらに2つの記事セットについてはその分類結果が被験者が与えた分類に含まれていた。分類が正しくなかった残り9記事セットのうち3つは、正解の分類がSingle-Productであるのに対してSingle-Eventと分類していた。実験に使用されたこれら記事セットの中に現れる固有物名（Product）は、特定の法案や国際条約等であり、記事セットに含まれる記事はその法案の審議や国際条約に対する発言について記述されたものであった。現在のアルゴリズムでは、Single-Eventを優先して分類するようになっているため、このような誤りが生じたと考えられる。しかしながら、正解の分類に関連する法案や国際条約等の固有物名は記事セット中の記事全体にわたって現れているため、判断手段101が一旦Single-Eventの分類を与えておきながらその後の判断過程を継続し、Single-Productの分類を与え直すことができるように構成することは可能であると考えられる。

【0037】

別の3記事セットでは、正解の分類がSingle-Eventであるのに対して異なる分類を付与していた。これらの記事セットでは、イベント名（Event）にあたる記述が固有表現ではなく、句や節の形で表されていた。一例を挙げると、「クリントン前大統領のホワイトハウス元実習生モニカ・ルインスキさんに対する不倫疑惑」という表現は固有表現ではないが、特定のイベントを指す表現である。現状の固有表現抽出システムでは、このような表現を一の固有表現として認識することはできない。また、記事セットのタイムスパンは設定した閾値 T_s よりも大きかった（上記例では、1年以上）ため、判断手段101は当該記事セットにSingle-Eventの分類を付与することができなかった。このことは、Single-Eventの分類を確実に付与するためには現在用いている手がかりの他に新たな手がかりを用

10

20

30

40

50

いる必要があるということを示唆している。

【0038】

また、実験の過程で、記事セットの中には本質的に一以上の分類を付与し得るものがあることが分かった。その理由の一つは、Single-Event、Multi-Event等、イベント名に基づく分類を付与すべき記事セットには他の固有表現クラスに基づく分類を付与可能なケースも多いことである。イベントの多くは、特定の人名や組織名、地名等に関連することがしばしばであり、イベントに関する記事を集めた記事セットに対しそのイベントに関連する（イベント名以外の）固有表現に焦点を絞ることが可能である。もう一つの理由として、イベントにおけるSingleとMultiとの区別が難しいことが挙げられる。あるイベントの中には、幾つかの小さなイベントが包含されることがある。例えば、「シドニーオリンピック」に関する記事セットは、一つのスポーツイベントを対象とするものとしてSingle-Eventの分類を付与し得るが、この記事セット中に複数の種目の結果を報じる記事が含まれているならば、それらに着目することでMulti-Eventの分類を付与することも可能である。イベントの単位をどのように認識するかは、被験者の観点到に依存する。

10

【0039】

因みに、上述の実験では、固有表現抽出手段による結果出力に人手による修正を加えて固有表現抽出タスクにおける誤りを排除している。本実施形態の文書セット分類装置を用いて機械的に分類を付与するにあたり、固有表現抽出タスクの段階でエラーが生起することを完全に避けるのは難しい（完璧な固有表現抽出ソフトウェアは現存しないため）。固有表現抽出手段における固有表現抽出処理を補完するためには、共参照や、文献（Y. Shin-yama, S. Sekine K. Sudo, and R. Grishman, "Automatic paraphrase acquisition from news articles", In Proceedings of the HLT-2002 conference, 2002.）に述べられているようなイベントの記述に関する言い換え表現の認識手法等を導入することが考えられる。

20

【0040】

以上では、入力として与えられる複数の記事（文書）が予め記事セット（文書セット）に仕分けされていることを前提としていた。しかしながら、入力として複数の記事が単純に与えられるような状況も考えられる。このような場合において、文書セット分類装置が、与えられる複数の記事を一または複数の記事セットに仕分けし、仕分けした記事セットに分類を付与するまでの処理を機械的に実行し得ることが好ましい。即ち、文書セット分類装置を構成するコンピュータ1が、図2に示す文書セット生成手段103としての機能をも発揮し得ることが好ましい。

30

【0041】

文書セット生成手段103は、ソフトウェアを主体として構成され、入力として与えられた文書の中に存在するキーワードを抽出し、一の文書のキーワードと他の文書のキーワードとの類似度を算出し、その類似度が閾値を超える場合にこれらの文書を同一の文書セットに割り当てることを通じて、複数の文書から少なくとも一の文書セットを生成する処理を行う。文書セット生成手段103が実行する処理の手順は、既に述べた記事セットの生成手法に類似する。即ち、プロセッサ1aが、プログラムに基づき、入力として与えられた記事データ（通常、メインメモリ1bまたは補助記憶デバイス1cの所要の記憶領域に格納されている）のうちの一つを選択的に読み込み、この記事データよりキーワード列を抽出する。キーワードの抽出は、形態素解析ソフトウェアを利用して行うことができる。例えば、記事データを形態素解析した結果より時相名詞、副詞的名詞を除いた名詞のうち頻度が所定値（例えば、2）以上のものをキーワードとして抽出する。入力として与えられた各記事データについて上記の方法でキーワードを抽出した後、プロセッサ1aが、一の記事データに係るキーワード列と他の記事データに係るキーワード列との間の類似度を算出する。プロセッサ1aが算出する類似度の指標としては、Dice's coefficient、Jaccard measure、cosine similarity等を採用することができる。その上で、類似度の指標が所定値（例えば、0.5）以上である複数の記事を類似記事として一の記事セットに含めることを通じて、記事セットの生成を行う。

40

50

【0042】

ところで、文書セット分類装置が出力する記事セットの分類についての情報を参照することで、当該記事セットに適した要約アルゴリズムを選択し得る。よって、文書セット分類装置が出力する分類についての情報を参照し、この分類に対応した要約アルゴリズムを用いて記事セット（文書セット）に含まれる複数記事（文書）の要約を行う要約装置を構築すれば、要約の精度の向上を図ることができる。

【0043】

本実施形態における文書要約装置は、文書セット分類装置を構成するコンピュータ1またはこのコンピュータ1とは別のコンピュータ（図示せず）に所定のプログラムをインストールすることで構築される。通常、プログラムは補助記憶デバイス1cに格納され、その実行の際には補助記憶デバイス1cからメインメモリ1bに読み込まれてプロセッサ1aにより解読される。そして、該プログラムに従いハードウェア資源を作動して、図3に示す要約手段201としての機能を発揮するようにしている。

【0044】

要約手段201は、文書セット分類装置が出力する、記事セットに付与された分類についての情報を参照し、この分類に対応した要約アルゴリズムにより前記記事セットに含まれる複数の記事を単一の文書に要約する。即ち、プロセッサ1aが、プログラムに基づき、記事セットに付与された分類についての情報を参照し、この分類に対応した要約アルゴリズムを選択する。しかる後、入力として与えられた記事セットに含まれる記事のデータを読み込み、要約の生成を行う。本実施形態における要約手段201は、ソフトウェアを主体として構成される。要約手段201の主体となるソフトウェアには、既知の複数記事要約ソフトウェアを応用できる。既知の複数記事要約ソフトウェアでは、一般に、記事セットに含まれる複数記事から重要と判断される文（ないし、文章）を抽出し、抽出した文を基にして要約を生成する。それぞれの文の重要度は、文の位置、文の長さ、文中に出現する単語の頻度、見出しとの類似度等の複数の条件に関するスコアを加算して算定される。また、重要度の算定にあたっては、個々の条件に関するスコアに対して重み付けがなされる（重みは、トレーニングデータを用いた訓練を通じて得られる）。しかして、本実施形態における要約手段201では、文書セット分類装置によって付与された分類に基づくスコアを加味して、それぞれの文の重要度を算出することとしている。具体例を挙げると、対象の記事セットに付与された分類がSingle-Personでありその主題を表す固有表現が「小泉」である場合には、各記事において、人名として認識できる「小泉」を含む文にスコアを与える。あるいは、対象の記事セットに付与された分類がMulti-Organizationである場合には、各記事において、組織名を含む文にスコアを与える。このように、分類に基づくスコアを加味してそれぞれの文の重要度を算定することにより、生成される要約の的確性の向上が期待できる。

【0045】

加えて、要約手段201が、複数記事中の一の文と他の文との間の類似度を（共通する単語の個数等を参照することで）算出して類似する複数の文を抽出し、抽出した類似する文のうち要約生成に用いる文を選出するものとしてもよい。その上で、類似する複数の文より要約生成に用いる文を選出するための処理を、対象の記事セットに付与された分類に応じて変更することが好適である。具体例を挙げると、対象の記事セットに付与された分類がSingle-classである場合には、類似する複数の文が同一の事物を表現している可能性が高いことから、類似する複数の文のうちの一部の文のみを代表として選出する。つまり、類似する文として、「京都市で起きた震度4の地震で、3人が軽いけがを負った。」、「京都市で起きた震度4の地震で、新たに2人が入院し、けが人は5人となった。」というような複数の文が抽出されたとき、これらのうち何れか一文のみを要約記事の要素として選出する。これらの文のうちの何れを選択するかは、それぞれの文の重要度のスコアを参照する、時系列で最も後者の文を選択する等のヒューリスティクスにより決定できる。他方、対象の記事セットに付与された分類がMulti-classである場合には、類似する複数の文が相異なる事物を表現している可能性が高いことから、重要度スコアの高い文に類似

10

20

30

40

50

する一部または全部の文をまとめて選出する。つまり、類似する文として、「京都市で起きた震度4の地震で、3人が軽いけがを負った。」、「大阪市で起きた震度5の地震で、5人が入院し、8人が軽いけがを負った。」というような複数の文が抽出されたとき、これらの文は表現上似ているものの相異なるイベントを記述していると考えられる。であるから、これらの文の全てを要約記事の要素として選出することもあり得る。

【0046】

総じて言えば、本実施形態における文書要約装置の要約手段201は、記事セットに付与された分類に基づく重要度スコアを加味してそれぞれの文の重要度を算定するプロセス、及び/または、記事セットに付与された分類に応じて類似する複数の文の取捨選択の手法を変えるプロセスを、既存の複数記事要約ソフトウェアに追加したものと構成可能である。そして、対象の記事セットに付与された分類に応じて異なる要約アルゴリズムの要約処理を実行可能である。因みに、Single-Personの分類が付与された記事セットより伝記的な記述を要約出力させたり、Multi-Productの分類が付与された記事セットより製品の名称・機能・値段等の要素を抽出させて表の形態で出力させたりというように、記事セットの分類に応じた多様な要約を要約手段201に出力させることも考えられる。

10

【0047】

本実施形態によれば、複数の文書の集合である文書セットに対し分類を付与するものとして、前記文書セットの主題が単独の固有表現(Named Entity)に関するものか複数の固有表現に関するものかを判断し、かつ、該固有表現が何れの固有表現クラスに属するかを判断する判断手段101と、前記判断手段101が下した判断に基づき、前記文書セットの主題に係る固有表現が単独であるか複数であるか、及び、該固有表現が属している固有表現クラスという2つの要素より規定される分類についての情報を出力する出力手段102とを具備する文書セット分類装置を構成したため、網羅性の高い分類を文書セットに付与可能となる。

20

【0048】

前記判断手段101が、前記文書セットに含まれる複数の文書の中に出現する固有表現の頻度または固有表現クラスの頻度のうち少なくとも一方を材料として、前記文書セットの主題に係る固有表現が単独であるか複数であるかの判断及び該固有表現が属する固有表現クラスの判断を実行するため、大文字で始まる語の頻度やhe、she等の人称代名詞の頻度等の特定言語の特性に依存することなく分類を実行し得る。即ち、より一般的に記事セットの分類を行うことができる。

30

【0049】

前記判断手段101が、前記文書セットに含まれる複数の文書の中に出現するクラスタムの頻度またはクラスタムが関連する固有表現クラスの頻度のうち少なくとも一方を材料として、前記文書セットの主題に係る固有表現が属する固有表現クラスの判断を実行するため、固有表現のみを手がかりとして分類できない文書セットに対しても適切な分類を付与することが可能である。

【0050】

また、前記判断手段101が、前記文書セットに含まれる複数の文書の各々の作成若しくは発表された時点に関する情報を参照し、これら複数の文書のうちの一部または全部が予め定められた期間内に作成若しくは発表されていることを条件として、前記記事セットの主題に係る固有表現が単独でありかつその属する固有表現クラスがイベント名クラスである旨の判断を下すものとしており、少なくともSingle-Eventクラスの文書セットを速やかに分類できる。

40

【0051】

文章セット分類装置が、与えられた文書の中に存在するキーワードを抽出し、一の文書のキーワードと他の文書のキーワードとの類似度を算出し、その類似度が閾値を超える場合にこれらの文書を同一の文書セットに割り当てることを通じて、複数の文書から少なくとも一の文書セットを生成し得る文書セット生成手段103をさらに具備するものであれば、与えられた複数の文書を一または複数の文書セットに仕分けしこれに分類を付与する

50

までの処理を一括に実行可能となる。このものは、与えられた文書を基に一または複数の要約を自動生成するシステムを構築するために有用となる。

【 0 0 5 2 】

さらに、上記の文書セット要約装置が文書セットに対して付与した分類を参照し、この分類に対応した要約アルゴリズムを選択して前記文書セットに含まれる複数の文書を要約する要約手段 2 0 1 を具備する文書要約装置を構成して、より適切な複数文書の要約処理を可能とすることができる。

【 0 0 5 3 】

なお、本発明は以上に詳述した実施形態に限られるものではない。特に、本発明で定義した記事セットの分類やその分類を行う文書セット分類装置及びそのプログラムは、自動要約以外での応用も考えられる。例示すると、情報検索において、検索された結果中の上位の記事を用いた検索後の再ランク付けや検索結果の効率的な表示を行うことができる。さらに、オープンドメインの情報抽出に利用することも考えられる。従来の情報抽出ではドメインが限定されており、記事の主題や分類は前提として与えられていた。しかし、ドメインを限定することなく情報抽出を行うためには、対象となるドメインの情報、即ち記事セットの分類を動的に実施する必要があると考えられるからである。

【 0 0 5 4 】

その他各部の具体的構成や図 1 0 ないし図 1 4 に示す処理の手順等もまた、上記実施形態に限られるものではなく、本発明の趣旨を逸脱しない範囲で種々変形が可能である。勿論、パーソナルコンピュータその他の汎用的な情報処理装置にプログラムをインストールして本発明に係る文書セット分類装置を構成することが可能であって、専用の装置を製造することが必須であるわけではない。

【 図面の簡単な説明 】

【 0 0 5 5 】

【 図 1 】 本発明に係る文書セット分類装置の構成説明図。

【 図 2 】 本発明に係る文書セット分類装置の構成説明図。

【 図 3 】 本発明に係る文書セット分類装置の構成説明図。

【 図 4 】 複数記事を包含する記事セットを例示する図。

【 図 5 】 2 人の被験者による記事分類実験の結果を示す図。

【 図 6 】 2 人の被験者による記事分類実験の結果を示す図。

【 図 7 】 本発明の一実施形態における文書セット分類装置が具備するハードウェア資源を示す図。

【 図 8 】 固有表現抽出手段による固有表現抽出処理を説明する図。

【 図 9 】 クラスタム抽出手段によるクラスタム抽出処理を説明する図。

【 図 1 0 】 判断手段が実行する判断処理の手順を示すフローチャート。

【 図 1 1 】 同フローチャート。

【 図 1 2 】 同フローチャート。

【 図 1 3 】 同フローチャート。

【 図 1 4 】 同フローチャート。

【 図 1 5 】 文書セット分類装置による自動分類実験の結果を示す図。

【 符号の説明 】

【 0 0 5 6 】

1 ... コンピュータ (文書セット分類装置、文書要約装置)

1 0 1 ... 判断手段

1 0 2 ... 出力手段

1 0 3 ... 文書セット生成手段

2 0 1 ... 要約手段

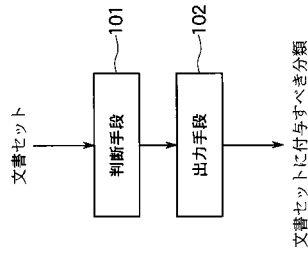
10

20

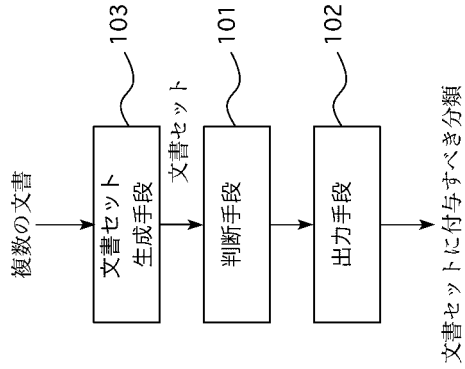
30

40

【 図 1 】



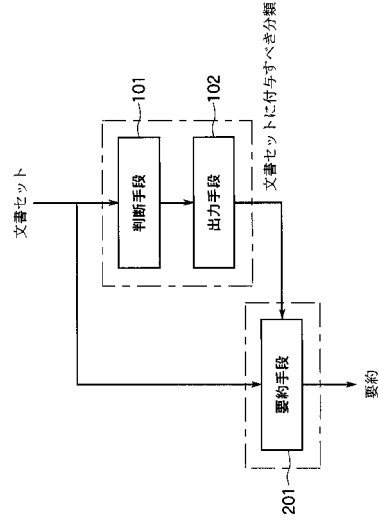
【 図 2 】



【 図 4 】

記事セットNo.	記事数	記事セットの主題	分類
1	14	ロシアのチェチェンとの紛争	Single-Event
2	10	日本各地の地震速報	Multi-Event
3	8	愛大薬学部失窃事件	Single-Event
4	7	民主リベラル新党の結成	Single-Event
5	6	書者初め大会、辯論大会等	Multi-Event
6	6	各国の首脳の動向	Multi-Person
7	6	日米首脳会談	Single-Event
8	4	大相撲の記事。特に、賞力花について	Single-Person
9	4	伊美議員の北海道知事選出馬	Single-Person
10	4	インドネシア・ポルトガル首脳会談	Single-Event
11	3	浄土真宗本願寺派の記事	Single-Event
12	3	市立船橋の高校サッカー大会での活躍	Single-Organization
13	3	ピエール・カー各首領のコメント	Multi-Organization
14	3	動物が関与した事件の記事	Other
15	3	レコード大賞、ノーベル賞に関する記事	Multi-Artifact
16	3	中央アルプス・木曽駒ヶ岳付近で発生した雪崩	Single-Event
17	3	関西国際空港に関する記事	Single-Location
18	3	サッカーのインターコンチネンタル選手権	Single-Event
19	3	野球選手のメジャーリーグ挑戦	Single-Person

【 図 3 】



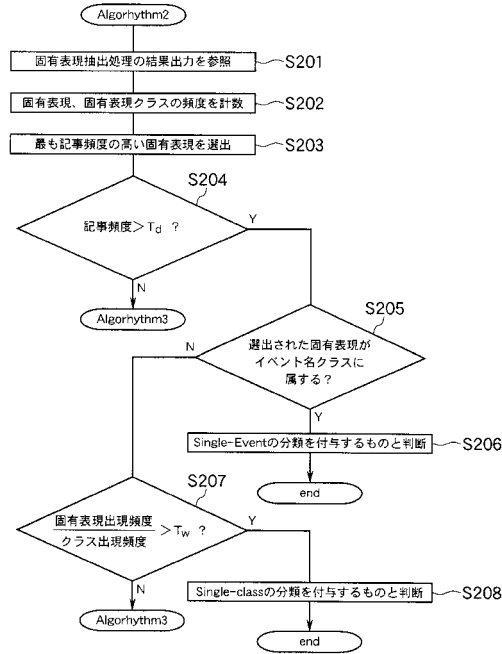
【 図 5 】

分類	被験者A	被験者B
Single-Person	2	1
Single-Organization	2	1
Single-Product	1	2
Single-Event	7	7
Multi-Location	1	1
Multi-Product	1	1
Multi-Event	6	9

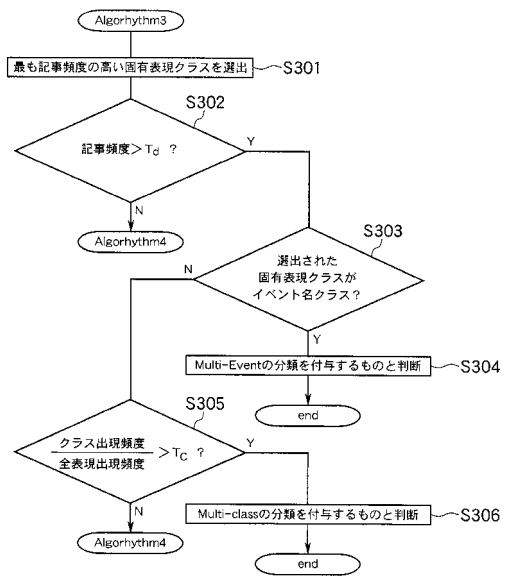
【 図 6 】

分類	被験者A	被験者B
Single-Person	12	10
Single-Location	2	4
Single-Product	3	5
Single-Event	6	6
Multi-Event	7	5

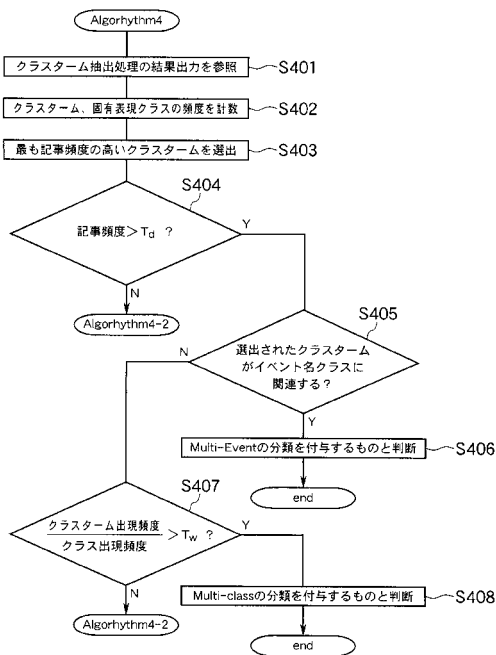
【 図 1 1 】



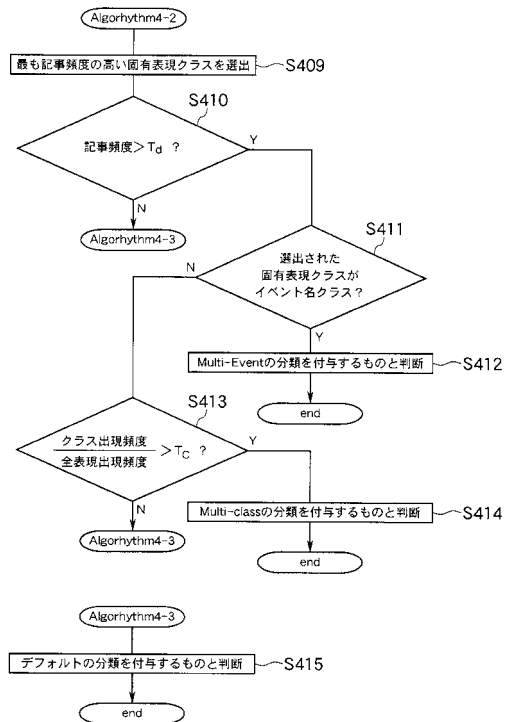
【 図 1 2 】



【 図 1 3 】



【 図 1 4 】



【 図 1 5 】

	一致	部分一致	誤り
文書セット分類装置	9(45%)	11(55%)	9
被験者A	14(70%)	17(85%)	3
被験者B	15(75%)	18(90%)	2
ベースライン	7(35%)	7(35%)	13

フロントページの続き

(72)発明者 関根 聡

アメリカ合衆国 10003 ニューヨーク州 ニューヨーク ブロードウェイ 715 セブンスフロア ニューヨーク大学内

審査官 丹治 彰

(56)参考文献 特開2001-331529(JP,A)

特開平04-106663(JP,A)

Kathleen R. McKeown et al., Columbia Multi-Document Summarization: Approach and Evaluation, Proceedings of DUC2001, 2001年 9月14日, [平成18年5月17日検索], URL, http://www-nlpir.nist.gov/projects/duc/pubs/2001papers/columbia_redo.pdf

難波英嗣ほか, テキスト自動要約 - 知的活動支援の基本技術として - 1 ここまで来たテキスト自動要約, 情報処理, 日本, 社団法人情報処理学会, 2002年12月15日, 第43巻 第12号, 第1287頁~第1294頁

(58)調査した分野(Int.Cl., DB名)

G06F 17/30

JSTPlus(JDream2)