

(19) 日本国特許庁(JP)

(12) 特 許 公 報(B2)

(11) 特許番号

特許第5234992号
(P5234992)

(45) 発行日 平成25年7月10日(2013.7.10)

(24) 登録日 平成25年4月5日(2013.4.5)

(51) Int. Cl. F 1
G 0 6 F 17/30 (2006.01)
 G 0 6 F 17/30 1 8 0 A
 G 0 6 F 17/30 2 1 0 D
 G 0 6 F 17/30 2 1 0 A

請求項の数 8 (全 21 頁)

(21) 出願番号	特願2009-121438 (P2009-121438)	(73) 特許権者	000004226
(22) 出願日	平成21年5月19日(2009.5.19)		日本電信電話株式会社
(65) 公開番号	特開2010-271800 (P2010-271800A)		東京都千代田区大手町二丁目3番1号
(43) 公開日	平成22年12月2日(2010.12.2)	(74) 代理人	100087446
審査請求日	平成23年10月12日(2011.10.12)		弁理士 川久保 新一
		(72) 発明者	川島 晴美
			東京都千代田区大手町二丁目3番1号 日 本電信電話株式会社内
		(72) 発明者	甲谷 優
			東京都千代田区大手町二丁目3番1号 日 本電信電話株式会社内
		審査官	野崎 大進

最終頁に続く

(54) 【発明の名称】 回答文書分類装置、回答文書分類方法及びプログラム

(57) 【特許請求の範囲】

【請求項1】

質問文書に対する回答文書を、グループを代表するキーワードである分類キーワードを含むグループ毎に分類する回答文書分類装置において、

単語と品詞との組み合わせによって構成されているパターンであって、所定の条件を示すパターンである条件パターンを蓄積している条件パターン蓄積手段と；

入力テキスト集合に含まれている文が、上記条件パターン蓄積手段に蓄積されている条件パターンを含んでいるか否かを判定する条件パターン判定手段と；

上記条件パターンを含んでいれば、上記条件パターンの前方に配置されている単語の配列である条件キーワードを抽出する条件キーワード抽出手段と；

上記条件パターンを含んでいなければ、主格になる文節に含まれている単語である主格キーワードを抽出する主格キーワード抽出手段と；

上記条件キーワードおよび上記主格キーワードに含まれるキーワード毎に、着目しているキーワードが上記条件キーワードとして出現する回数である出現回数と、上記主格キーワードとして出現する回数である出現回数とを集計し、上記条件キーワードの出現回数と上記主格キーワードの出現回数とが所定の条件を満たす場合に、上記着目しているキーワードを分類キーワードであると判定する分類キーワード判定手段と；

を有することを特徴とする回答文書分類装置。

【請求項2】

質問文書に対する回答文書を、グループを代表するキーワードである分類キーワードを

含むグループ毎に分類する回答文書分類装置において、

単語と品詞との組み合わせによって構成されているパターンであって、所定の条件を示すパターンである条件パターンを蓄積している条件パターン蓄積手段と；

回答文書集合に含まれている文が、上記条件パターン蓄積手段に蓄積されている条件パターンを含んでいるか否かを判定する条件パターン判定手段と；

上記条件パターンを含んでいれば、上記条件パターンの前方に配置されている単語の配列である条件キーワードを抽出する条件キーワード抽出手段と；

上記条件パターンを含んでいなければ、主格になる文節に含まれている単語である主格キーワードを抽出する主格キーワード抽出手段と；

回答文書の文章毎に、条件キーワードであることを示す情報または主格キーワードであることを示す情報が追加されている解析結果と、回答文書を一意に示す回答文章IDとを対応付けた分類キーワード候補を蓄積する分類キーワード候補蓄積手段と；

回答文章IDが複数与えられた際に、回答文章IDに該当する回答文書の解析結果の集合を、上記分類キーワード候補蓄積手段から取得し、上記解析結果の集合に含まれている条件キーワードおよび主格キーワードに含まれるキーワード毎に、着目しているキーワードが上記条件キーワードとして出現する回数である出現回数と、上記主格キーワードとして出現する回数である出現回数とを集計し、上記条件キーワードの出現回数と上記主格キーワードの出現回数とが所定の条件を満たす場合に、上記着目しているキーワードを分類キーワードであると判定する分類キーワード判定手段と；

を有することを特徴とする回答文書分類装置。

【請求項3】

請求項1または請求項2において、

上記条件キーワード抽出手段は、上記条件パターンを含む文節から、条件パターンを除いた語句に、条件パターンを含む文節へ係る0個以上の文節を追加した語句を、条件キーワードとして抽出する手段であることを特徴とする回答文書分類装置。

【請求項4】

請求項1または請求項2において、

上記主格キーワード抽出手段は、主格になる文節に含まれている名詞句に、主格になる文節へ係る0個以上の文節を追加した語句を、上記条件キーワードとして抽出する手段であることを特徴とする回答文書分類装置。

【請求項5】

請求項1～請求項4のいずれか1項において、

上記分類キーワード判定手段は、上記条件キーワードとしての出現回数が多いほど、上記分類キーワードとして判定され易くなる条件を用いて、上記分類キーワードを判定する手段であることを特徴とする回答文書分類装置。

【請求項6】

質問文書に対する回答文書を、グループを代表するキーワードである分類キーワードを含むグループ毎に分類する回答文書分類方法において、

単語と品詞との組み合わせによって構成されているパターンであって、所定の条件を示すパターンである条件パターンが蓄積されている条件パターン蓄積手段に蓄積されている条件パターンを、入力テキスト集合に含まれている文が含まれているか否かを、条件パターン判定手段が判定し、記憶装置に記憶する条件パターン判定工程と；

上記条件パターンを含んでいれば、上記条件パターンの前方に配置されている単語の配列である条件キーワードを、条件キーワード抽出手段が抽出し、記憶装置に記憶する条件キーワード抽出工程と；

上記条件パターンを含んでいなければ、主格になる文節に含まれている単語である主格キーワードを、主格キーワード抽出手段が抽出し、記憶装置に記憶する主格キーワード抽出工程と；

上記条件キーワードおよび上記主格キーワードに含まれるキーワード毎に、着目しているキーワードが上記条件キーワードとして出現する回数である出現回数と、上記主格キー

10

20

30

40

50

ワードとして出現する回数である出現回数とを集計し、上記条件キーワードの出現回数と上記主格キーワードの出現回数とが所定の条件を満たす場合に、上記着目しているキーワードを分類キーワードであると、分類キーワード判定手段が判定し、記憶装置に記憶する分類キーワード判定工程と；

を有することを特徴とする回答文書分類方法。

【請求項 7】

質問文書に対する回答文書を、グループを代表するキーワードである分類キーワードを含むグループ毎に分類する回答文書分類方法において、

単語と品詞との組み合わせによって構成されているパターンであって、所定の条件を示すパターンである条件パターンを蓄積している条件パターン蓄積手段に蓄積されている条件パターンを、回答文書集合に含まれている文が含んでいるか否かを、条件パターン判定手段が判定し、記憶装置に記憶する条件パターン判定工程と；

10

上記条件パターンを含んでいれば、上記条件パターンの前方に配置されている単語の配列である条件キーワードを、条件キーワード抽出手段が抽出し、記憶装置に記憶する条件キーワード抽出工程と；

上記条件パターンを含んでいなければ、主格になる文節に含まれている単語である主格キーワードを、主格キーワード抽出手段が抽出し、記憶装置に記憶する主格キーワード抽出工程と；

回答文書の文章毎に、条件キーワードであることを示す情報または主格キーワードであることを示す情報が追加されている解析結果と、回答文書を一意に示す回答文章IDとを対応付けた分類キーワード候補とを、分類キーワード候補蓄積手段が蓄積し、記憶装置に記憶する分類キーワード候補蓄積工程と；

20

回答文章IDが複数与えられた際に、回答文章IDに該当する回答文書の上記解析結果の集合を、上記分類キーワード候補蓄積工程から取得し、上記解析結果の集合に含まれている条件キーワードおよび主格キーワードに含まれるキーワード毎に、着目しているキーワードが上記条件キーワードとして出現する回数である出現回数と、上記主格キーワードとして出現する回数である出現回数とを集計し、上記条件キーワードの出現回数と上記主格キーワードの出現回数とが所定の条件を満たす場合に、上記着目しているキーワードを分類キーワードであると、分類キーワード判定が判定し、記憶装置に記憶する分類キーワード判定工程と；

30

を有することを特徴とする回答文書分類方法。

【請求項 8】

請求項 1～請求項 5 のいずれか 1 項に記載の回答文書分類装置を構成する各手段としてコンピュータを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、質問についての回答を分類する装置に関し、特に、質問についての複数の回答を、キーワード毎に分類して提供する回答文書分類装置に関する。

【背景技術】

40

【0002】

従来、インターネット等のネットワーク上で公開されている情報を探す手段として、キーワード検索がよく利用されている。キーワード検索は、利用者から 1 つまたは複数の検索キーワードを入力し、入力された検索キーワードの全て、またはいずれかを含む文書を検索して出力するサービスである。検索結果として出力される文書の数は、検索キーワードが一般的な語であるほど数が多く、利用者が目的に合った情報を見つけるのは困難である。そのために、利用者は検索キーワードを追加したり、調べたい分野に特徴的な検索キーワードを指定したりして、検索結果を絞り込む作業を行なっている。

【0003】

また、キーワード検索でうまく目的に合った情報を見つけれない場合や、人に相談し

50

たいような場合に、利用者同士が質問と回答をやり取りするQ & Aコミュニティーサイトを利用する場合がある。Q & Aコミュニティーサイトでは、図5に示すように、質問したい利用者が投稿した1つの質問32に対して、複数の回答者によって、回答33、34、35が投稿され、サイト上でコミュニケーションが行われる。Q & Aコミュニティーサイトにおいてもキーワード検索する機能が提供され、検索する対象を質問だけ、質問・回答の両方等と指定して検索することができる。この場合も、検索キーワードが一般的な語であれば検索結果の数は多くなり、目的に合った情報を探すのは容易ではない。

【0004】

目的に合った情報を探しやすくするために、検索結果をグループ化して提示する手法がある。たとえば、利用者が入力した検索キーワードに対して、過去に入力された検索キーワードの履歴の中から、利用者が入力した検索キーワードに隣接して良く入力されるキーワードを関連語として抽出し、該関連語毎にグループ化して検索結果を表示する発明が知られている(たとえば、特許文献1参照)。この発明によれば、利用者が自分で検索結果を絞り込むためのキーワードを考える必要がなく、選択するだけで良いという利点がある。

10

【0005】

また、回答文書を分類する方法として、クラスタリング技術を利用することができる。クラスタリング技術は、文章集合が与えられると、文章間の類似度を算出し、類似した文章毎にグループを作成する方法である。文章間の類似度は、単語出現頻度に基づく文章ベクトルで文章を表し、文章ベクトル間のコサイン類似度を適用する手法が広く用いられている。すなわち、文章 d_n を文章ベクトル

20

【0006】

【数1】

$$\vec{X}_n = (x_{n1}, x_{n2}, \dots, x_{nv})$$

【0007】

により表す場合、 v は、単語集合 $W = \{w_1, w_2, \dots, w_v\}$ 中の単語の総数を示し、 x_{ni} は、文章 d_n における単語 w_i の重みを示す。このとき、文章 d_j と文章 d_k の類似度は、各文章ベクトルがなす角

30

【0008】

【数2】

$$\cos \theta_{j,k} = \frac{\vec{X}_j \cdot \vec{X}_k}{|\vec{X}_j| |\vec{X}_k|}$$

【0009】

で表される。また w_i の重みは、単語の文章内での出現頻度 tf (term frequency)をそのまま利用する場合や、出現頻度 tf に、単語出現数を全文章数で割った値の対数 idf を乗算した $tf-idf$ (term frequency/inverse document)を利用する。つまり、類似度の高い文章同士は、この単語の重みの傾向が似通っている文章同士であることを意味する。

40

【先行技術文献】

【特許文献】

【0010】

【特許文献1】特許第4009937号公報

【発明の概要】

【発明が解決しようとする課題】

【0011】

50

本来、Q & A コミュニティーサイトの利用者の目的は、質問についての回答を得ることであるので、回答文書を分類して提供することが望ましい。上記特許文献 1 記載の従来技術を、Q & A コミュニティーサイトに適用する場合、検索する利用者は知りたい情報の分野に詳しくない場合が多く、適切なキーワードを入力しているとは限らない。このために、検索キーワードを分類用のキーワードとして利用するだけでは、回答文書を精度良く分類できないという問題がある。

【0012】

また、クラスタリング技術によって回答文書を分類する場合、文章数の二乗に比例して処理時間がかかるので、キーワード検索結果の文章集合が多い場合、分類結果が出力されるまでに時間がかかるという問題がある。

10

【0013】

本発明は、検索キーワードに含まれないキーワードも利用して、短時間で回答文書を分類することができる回答文書分類装置、回答文書分類方法および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】

【0014】

本発明は、質問文書に対する回答文書を、グループを代表するキーワードである分類キーワードを含むグループ毎に分類する回答文書分類装置において、単語と品詞との組み合わせによって構成されているパターンであって、所定の条件を示すパターンである条件パターンを蓄積している条件パターン蓄積手段と、入力テキスト集合に含まれている文が、上記条件パターン蓄積手段に蓄積されている条件パターンを含んでいるか否かを判定する条件パターン判定手段と、上記条件パターンを含んでいれば、上記条件パターンの前方に配置されている単語の配列である条件キーワードを抽出する条件キーワード抽出手段と、上記条件パターンを含んでいなければ、主格になる文節に含まれている単語である主格キーワードを抽出する主格キーワード抽出手段と、上記条件キーワードおよび上記主格キーワードに含まれるキーワード毎に、着目しているキーワードが上記条件キーワードとして出現する回数である出現回数と、上記主格キーワードとして出現する回数である出現回数とを集計し、上記条件キーワードの出現回数と上記主格キーワードの出現回数とが所定の条件を満たす場合に、上記着目しているキーワードを分類キーワードであると判定する分類キーワード判定手段とを有する回答文書分類装置である。

20

30

【発明の効果】

【0015】

本発明によれば、質問についての回答文書を検索したい利用者が、検索キーワードとして入力しないキーワードも利用して、回答文書を分類することができ、質問者は、分類された情報を選択するだけで、所望の情報を取得することができるという効果を奏する。

【0016】

また、本発明によれば、回答文書集合に含まれるキーワードの出現回数に基づいて分類キーワードの判定を行うので、処理が高速であるという効果を奏する。

【図面の簡単な説明】

【0017】

【図1】本発明の原理を説明する図である。

【図2】本発明の実施例 1 である回答文書分類装置 100 の構成を示す図である。

【図3】条件パターン蓄積手段 11 が蓄積している条件パターンの例を示す図である。

【図4】条件パターン判定手段 12 の処理を示すフローチャートである。

【図5】Q & A コミュニティーサイトの構成を示す図である。

【図6】条件パターン判定手段 12 が判定した形態素解析結果を示す図である。

【図7】実施例 1 において、条件キーワード抽出手段 13 が抽出した係り受け解析した結果の例を示す図である。

【図8】主格キーワード抽出手段 14 の動作を示す図である。

【図9】分類キーワード判定手段 15 の動作の説明図である。

40

50

【図10】分類キーワード判定手段19の動作を示すフローチャートである。

【図11】本発明の実施例2である回答文書分類装置200を示すブロック図である。

【図12】実施例2における分類キーワード蓄積手段18の蓄積例を示す図である。

【発明を実施するための形態】

【0018】

発明を実施するための形態は、以下の実施例である。

【実施例1】

【0019】

本発明では、グループを代表するキーワードであり、回答文書を分類するためのキーワードである分類キーワードを、回答文書の中から抽出する。質問された内容に詳しい回答者は、質問文書中に含まれていない語句を用いて、詳細な情報を提供することが考えられるので、上記分類キーワードを回答文書から抽出する。質問文書についての回答文書を、分類キーワード毎にグループ化して利用者に提供する。これによって、利用者は、分類キーワードを手がかりに、所望の情報を容易に取得することができる。

10

【0020】

本発明は、質問文書についての回答文書を、分類キーワード毎にグループ化する回答文書分類装置において、与えられた回答文書集合に含まれているキーワード(単語)から、グループを代表するキーワードである分類キーワードを判定する。

【0021】

まず、回答文書の特徴について説明する。たとえば、次のような質問文書があった場合

20

【0022】

「6月にテーマパークXに遊びに行こうと計画しています。雨に備えて傘か合羽を用意しようと思っていますが、どちらがいいでしょう？」

この質問に対して3人の回答者がそれぞれ以下の回答をしたとする。

【0023】

回答者A：「傘は人が多いと迷惑になるので、合羽がお勧めです。」

回答者B：「大人だけなら、傘で大丈夫ですよ。小さい子供が一緒なら、合羽が楽です。」

30

回答者C：「大人は傘が良いと思うよ。」

回答者Aは、質問文書に記載された内容の範囲で自分の意見を述べている。回答者Bは、質問文書には記載されていない「大人だけ」なのか「子供が一緒」なのかの情報を追加して、それぞれの場合について自分の意見を述べている。

【0024】

質問者は、「大人だけ」か「子供が一緒」かによって、お勧めの情報が異なることを、質問した時点では知らないが、回答文書の分類キーワードとして、「大人だけ」、「子供が一緒」というキーワードが提示されれば、自分の状況に合わせて分類キーワードを選択することができ、的確な回答文書を参照することが可能になる。

【0025】

回答者Cは「大人は」と述べ、条件パターンである「 なら」を用いていない。なお、上記「条件パターン」は、明示的に条件であることが分かる文節である。

40

【0026】

回答者Cによる回答文書中、「大人は」という主格となる文節に、暗黙的に「大人の場合」という条件を含めて記述している。回答文書において、八格である「 は」は、 について説明をする場合に用いられる場合が多い。そこで、「 は」という文節の中にも、回答文書を分類する場合に適したキーワードが含まれていると考える。

【0027】

図1は、本発明の原理を説明するフローチャートである。

【0028】

まず、入力されたテキスト集合に含まれている各文に対して、上記条件パターンが記述

50

されているか否かを判定する（S1）。上記「条件パターン」は、たとえば「 なら、
××です。」という文における「 なら」中の「なら」である。

【0029】

なお、条件パターンとして、次のパターンもが考えられる。

【0030】

「 であれば」中の「であれば」、
「 の場合」中の「の場合」、
「 ですと」中の「ですと」。

【0031】

文中に、上記条件パターンが記述されていると判定されると、「 なら」という条件を示す文節中の単語「 」に含まれているキーワードを、「条件キーワード」として抽出する（S2）。つまり、「条件キーワード」は、条件パターンを含む文節から、条件パターンを除いた単語に含まれているキーワードである。また、条件パターンが記述されていないと判定されると、「 は」のように、主格になる文節が存在するかどうかを調べ、主格になる文節が存在すれば、主格になる文節「 」に含まれているキーワードを、「主格キーワード」として抽出する（S3）。

10

【0032】

つまり、「主格キーワード」は、主格になる文節に含まれている名詞句である。

【0033】

条件キーワードと主格キーワードとは、少なくとも1つの単語によって構成されている

20

【0034】

次に、与えられている入力テキスト集合全体から抽出したキーワードのうちで、着目しているキーワードが、上記条件キーワードとして出現する回数である出現回数と、上記主格キーワードとして出現する回数である出現回数とを求め、上記条件キーワードの出現回数と上記主格キーワードの出現回数とが所定の条件を満たす場合に、上記着目しているキーワードを分類キーワードであると判定する（S4）。

【0035】

つまり、「分類キーワード」は、条件キーワードの出現回数と、主格キーワードの出現回数とが、所定の条件を満たす場合におけるキーワードである。

30

【0036】

図2は、本発明の実施例1である回答文書分類装置100の構成を示す図である。

【0037】

回答文書分類装置100は、テキスト入力手段20とテキスト出力手段40とに、接続され、条件パターン蓄積手段11と、条件パターン判定手段12と、条件キーワード抽出手段13と、主格キーワード抽出手段14と、分類キーワード判定手段15とを有する。

【0038】

条件パターン判定手段12は、テキスト入力手段20から回答文書の集合を入力すると、入力した回答文書を文単位に分割し、各文に条件パターンが含まれているかどうかを判定し、そして、条件パターンが含まれている文を、条件キーワード抽出手段13へ送り、条件パターンが含まれていない文を、主格キーワード抽出手段14へ送る。

40

【0039】

条件キーワード抽出手段13は、条件パターン判定手段12から、条件パターンを含む形態素解析結果の文と、検出された条件パターンを示す情報とを受け取り、条件パターンに一致する箇所の前方に存在するキーワードを、条件キーワードとして抽出し、記憶装置に記憶する手段である。

【0040】

主格キーワード抽出手段14は、条件パターン判定手段12から、条件パターンを含まないと判定された文を受け取り、係り受け解析し、記憶装置に記憶する。この解析結果中に、主格となる文節が存在するかどうかを調べ、主格となる文節が存在すれば、主格とな

50

る文節に含まれているキーワード（名詞句）を、「主格キーワード」として抽出し、記憶装置に記憶する。

【0041】

分類キーワード判定手段15は、条件キーワード抽出手段13が抽出した条件キーワードと、主格キーワード抽出手段14が抽出した主格キーワードとを入力し、キーワード毎に、所定の条件（後述の式（1）、式（2）に示す条件）を満たすかどうかを判定し、この予め設定された条件を満たすキーワードを、分類キーワードとして、記憶装置に記憶し、この分類キーワードを出力する。分類キーワード判定手段15は、具体的には、分類キーワードと、このキーワードを含む回答文章IDと、文番号の情報とを組みにして、テキスト出力手段40に送る。

10

【0042】

なお、テキスト入力手段20とテキスト出力手段40との間に、質問・回答文書蓄積手段30が接続されている。質問・回答文書蓄積手段30は、Q & A コミュニティーサイトにおいて投稿された質問文書、回答文書を多数、蓄積している。

【0043】

テキスト入力手段20は、質問・回答文書蓄積手段30から、後述する特定の条件に合致する回答文書の集合を取得し、条件パターン判定手段12に送る。

【0044】

後述する特定の条件に合致する回答文書の集合を取得するための条件は、特定のキーワードを含む質問に回答した回答文書の集合等である。

20

【0045】

テキスト出力手段40は、分類キーワード判定手段15が出力した分類キーワードを入力し、分類キーワード毎に、回答文書を分類し、記憶装置に記憶し、出力する。

【0046】

図3は、条件パターン蓄積手段11が蓄積している条件パターンの例を示す図である。

【0047】

条件パターン蓄積部11は、条件パターンとして抽出したい単語と、この単語の出現順序パターンとを、図3に示すように予め蓄積する。図3に示す蓄積例では、1行に1つの条件パターンが記述されている。1つの条件パターンは、「品詞：読み」がカンマで接続され、左から順に連続して単語が出現するパターンである。4行目の単語パターン、「判定詞：デ、動詞語幹：ア、動詞活用語尾：レ、動詞接尾辞：バ」は、「であれば」という条件パターンを抽出するための記述である。

30

【0048】

条件パターン判定手段12は、テキスト入力手段20から回答文書の集合を入力すると、入力した回答文書を文単位に分割し、各文に条件パターンが含まれているかどうかを判定する。そして、条件パターンが含まれている文を、条件キーワード抽出手段13へ送り、条件パターンが含まれていない文を、主格キーワード抽出手段14へ送る。なお、各回答文書には、回答文書を一意に特定する回答文章IDがそれぞれ付与される。

【0049】

図4は、条件パターン判定手段12の動作を示すフローチャートである。

40

【0050】

まず、条件パターン蓄積手段11に蓄積されている条件パターンを全て読み込み（S51）、テキスト入力手段20から受け取った回答文書について、回答文書毎に、回答文章IDに対応する回答文書を文単位に分割し、この分割された文に、回答文書内での出現番号を付与する（S52、S53）。

【0051】

続いて、未処理の文番号が存在すれば（S54）、未処理の文番号を1つ選択し、形態素解析処理を行う（S55）。形態素解析結果の単語の品詞と読みとが、条件パターンに一致する箇所があるかどうかを調べ（S56）、条件パターンに一致する箇所があれば（S56、YES）、処理対象の回答文章IDと、文番号と、文と条件パターンとが一致し

50

た箇所とを、条件キーワード抽出手段13へ送る(S57)。条件パターンと一致する箇所がなければ(S56、NO)、処理対象の回答文章IDと、文番号と、文とを、主格キーワード抽出手段14へ送る(S58)。

【0052】

図5は、Q&Aコミュニティサイトの構成例を示す図である。

【0053】

図5に示す回答文書33、34、35を、条件パターン判定手段12が入力した場合における具体的な動作について説明する。回答文書33、34、35の文章IDを、それぞれ、33、34、35とする。まず、未処理の回答文章ID33が選択され(S52、S53)、この選択された回答文章ID33は、1文から構成されているので(S53)、1文が選択され、形態素解析処理が行われる(S54、S55)。

10

【0054】

図6は、条件パターン判定手段12が判定した形態素解析結果を示す図である。

【0055】

形態素解析結果61において、1行目に回答文章IDが記載され、2行目に文番号が記載され、3行目以降に、形態素毎の表記、品詞、読みが記載されている。形態素解析結果61の中に、条件パターンと一致する単語出現パターン(単語が出現するパターン)があるかどうかを調べる(S56)。図6に示す形態素解析結果61には、条件パターンと一致する箇所がないので(S56、NO)、形態素解析結果61を主格キーワード抽出手段14へ送る。

20

【0056】

回答文章ID33には他に文がないので(S54、NO)、次の未処理の回答文書34を選択する(S52)。回答文章ID34の文章は、2つの文に分割され(S53)、1番の文が選択され、形態素解析処理される(S55)。形態素解析結果62に、条件パターンに一致する単語出現パターンがあるかどうかを調べると(S56)、条件パターン「判定詞：ナラ」と一致する箇所63が発見される。

【0057】

条件パターンと一致する箇所63があるので(S56、YES)、形態素解析結果62中で、条件パターンと一致する箇所を明示するために、「*」を付与し、条件キーワード抽出部13へ送る(S57)。続いて、回答文章ID34の文章の2番目の文(形態素解析結果64)について、S54~S56の処理を実行し、条件パターンに一致する箇所65が発見される。

30

【0058】

次に、回答文章ID35の文章について、S52~S56の処理を実行し、形態素解析結果66に、条件パターンが存在しないので(S56、NO)、形態素解析結果66を、主格キーワード抽出手段14へ送る。未処理の回答文章IDがなくなったので(S52、NO)、処理を終了する。

【0059】

条件キーワード抽出手段13は、条件パターン判定手段12から、条件パターンを含む形態素解析結果の文と、条件パターンに一致する箇所(検出された条件パターン)を示す情報とを受け取り、条件パターンに一致する箇所の前方に存在するキーワードを、条件キーワードとして抽出する。

40

【0060】

条件キーワードを抽出する方法として、次の[方法1]~[方法3]が考えられる。

[方法1] 検出された条件パターンに最も近い前方の名詞を、条件キーワードとして抽出する方法、

[方法2] 検出された条件パターンを含む文節から、条件パターンを除いた語句を、条件キーワードとして抽出する方法、

[方法3] 上記[方法2]において、検出された条件パターンを含む文節に係る文節を、N個追加したものを、条件キーワードとして抽出する方法。

50

【 0 0 6 1 】

上記〔方法 3〕における整数値 N を、予め設定するようにしてもよく、利用者が設定できるようにしてもよい。

【 0 0 6 2 】

図 7 は、実施例 1 において、条件キーワード抽出手段 1 3 が抽出した係り受け解析した結果の例を示す図である。

【 0 0 6 3 】

つまり、図 7 は、回答文章 ID 3 4 の文章の文番号 1 の形態素解析結果 6 2 を、係り受け解析した結果 7 1 と、文番号 2 の形態素解析結果 6 4 を、係り受け解析した結果 7 2 とを示す図である。係り受け解析した結果 7 1 において、上記〔方法 1〕では、「大人」が、条件キーワードとして抽出され、上記〔方法 2〕では、「大人だけ」が、条件キーワードとして抽出され、上記〔方法 3〕では、「大人だけなら」に係る文節が存在しないので、「大人だけ」が、条件キーワードとして抽出される。

【 0 0 6 4 】

係り受け解析した結果 7 2 において、上記〔方法 1〕では、「一緒」が、条件キーワードとして抽出され、上記〔方法 2〕でも、「一緒」が、条件キーワードとして抽出され、上記〔方法 3〕では、条件パターンを含む文節に係る文節の数 $N = 1$ であれば、「子供が一緒」が、条件キーワードとして抽出される。抽出された条件キーワードは、抽出された回答文章 ID と文番号とに対応付けられ、分類キーワード判定手段 1 5 へ渡される。

【 0 0 6 5 】

主格キーワード抽出手段 1 4 は、条件パターン判定手段 1 2 から、条件パターンを含まないと判定された文を受け取り、係り受け解析を行う。この解析結果中に、主格となる文節が存在するかどうかを調べ、主格となる文節が存在すれば、主格となる文節に含まれているキーワード（名詞句）を、「主格キーワード」として抽出する。本発明では、主格となる文節として、八格「〇〇は」を抽出する。八格「は」は、動作主体となる場合もあるが、名詞句「」について説明する場合にも用いられる。このために、回答文書の中で、名詞句「」の説明をしている箇所は、質問者にとって有益な情報となる可能性があると考え、八格「は」を抽出することとする。

【 0 0 6 6 】

主格キーワードとして、八格の名詞句部分「」だけを抽出する場合と、条件キーワードを抽出する場合における上記〔方法 3〕のように、主格となる文節に係る文節を、N 個追加したものを、主格キーワードとして抽出するようにしてもよい。

【 0 0 6 7 】

図 8 は、主格キーワード抽出手段 1 4 の動作を示す図である。

【 0 0 6 8 】

図 8 (1) は、たとえば、回答文章 ID 3 3、3 5 の文が、主格キーワード抽出部 1 4 に入力された場合、回答文章 ID 3 3 を、係り受け解析した結果を示す図である。図 8 (2) は、回答文章 ID 3 5 の文を、係り受け解析した結果を示す図である。

【 0 0 6 9 】

図 8 (1) に示す例において、主格となる文節は、「傘は」であり、この文節から、名詞「傘」を、主格キーワードとして抽出する。これと同様に、図 8 (2) に示す例において、主格となる文節は、「大人は」であり、この文節から名詞「大人」を、主格キーワードとして抽出する。抽出された主格キーワード「傘」、「大人」を、それぞれ回答文章 ID と文番号とを対応付け、「分類キーワード」として分類キーワード判定手段 1 5 へ送る。なお、分類キーワードと回答文章 ID と文番号を対応づけて分類キーワード判定手段 1 5 へ送るのであり、回答文章 ID と文番号が分類キーワードに含まれるわけではない。

【 0 0 7 0 】

分類キーワード判定手段 1 5 は、条件キーワード抽出手段 1 3 が抽出した条件キーワードと回答文章 ID と文番号、主格キーワード抽出手段 1 4 が抽出した主格キーワードと回答文章 ID と文番号を入力し、キーワード毎に、後述の式 (1)、式 (2) に示す条件を

10

20

30

40

50

満たすかどうかを判定し、この予め設定された条件を満たすキーワードを、分類キーワードとして出力する。

【0071】

つまり、上記「分類キーワード」は、条件キーワードと、主格キーワードとが、予め設定された条件を満たすキーワードである。

【0072】

上記予め設定された条件は、キーワード w が条件キーワードとして出現する回数である出現回数 $X(w)$ と、キーワード w が主格キーワードとして出現する回数である出現回数 $Y(w)$ とに応じた条件である。たとえば、次の式(1)を条件とし、判定式 $F(w)$ を求める。

【0073】

$$F(w) = X(w) + (1 - \alpha) Y(w) \dots \text{式(1)}$$

$$F(w) \geq t_h \dots \text{式(2)}$$

ここで、 $0 < \alpha < 1$ であり、 α を 1 に近い値に設定すると、上記式(1)は、条件キーワードとしての出現回数が多いキーワードほど、高い値になる。さらに、上記式(2)によって、予め設定された閾値 t_h 以上の値をとるキーワードのみを、分類キーワードとして判定することができる。

【0074】

出現回数を集計するときに、厳密に一致するキーワードだけを、出現回数としてカウントする場合と、一部が一致する場合(部分一致)も、出現回数としてカウントする場合とが考えられる。以下では、部分一致した場合に出現回数としてカウントする場合について説明する。

【0075】

図9は、分類キーワード判定手段15の動作の説明図である。

【0076】

たとえば、条件キーワード抽出手段において「大人だけ」「小さな子供」の2語が抽出され、主格キーワード抽出手段14において「傘」「大人」の2語が主格キーワードとして抽出された場合について考える。

【0077】

条件キーワード情報91は、1行に「条件キーワード 回答文章ID 文番号」の順にスペース区切りで記載されている。主格キーワード情報92は、1行に「主格キーワード 回答文章ID 文番号」の順にスペース区切りで記載されている。

【0078】

図10は、分類キーワード判定手段15の動作を示すフローチャートである。

【0079】

分類キーワード判定手段15は、分類キーワード候補を選ぶ(S101)。分類キーワード候補は、入力された条件キーワードと主格キーワードとの中から、重複を除去したキーワードの集合である。

【0080】

集計結果93における分類キーワード候補の列に記載する4種類が選ばれる。次に、未処理の分類キーワード候補が存在すれば(S102、YES)、分類キーワード候補を1つ選択し、この分類キーワード候補の条件キーワード出現回数と主格キーワード出現回数とを0に初期化する(S103)。

【0081】

次に、入力された条件キーワードを、1つずつ「比較キーワード」として選択し、この選択された比較キーワードを、分類キーワード候補と比較する。なお、分類キーワード候補と比較する条件キーワードを、便宜上、「比較キーワード」と表現する。

【0082】

未処理の条件キーワードが存在すれば(S104)、未処理の条件キーワードを、比較キーワードとして、1つ選択し、分類キーワード候補と比較する(S105)。

10

20

30

40

50

【0083】

この比較の結果、完全一致の場合、または分類キーワード候補に比較キーワードを含む場合（部分一致の場合）には（S106、YES）、条件キーワード出現回数を1増やす（S107）。これ以外の場合は（S106、NO）、S104に戻り、未処理の条件キーワードがなくなるまで、S104～S107の処理を実行する。未処理の条件キーワードがなくなれば（S104、NO）、主格キーワード似ついで比較する処理に移行する。

【0084】

未処理の主格キーワードが存在すれば（S108、YES）、未処理の主格キーワードを、比較キーワードとして、1つ選択し、分類キーワード候補と比較する（S109）。この比較の結果、完全一致である場合、または分類キーワード候補に比較キーワードを含む場合（部分一致の場合）には（S110、YES）、主格キーワード出現回数を、1増やす（S111）。これ以外の場合（S110、NO）、S108に戻り、未処理の主格キーワードがなくなるまで、S108～S111の処理を実行する。

10

【0085】

未処理の主格キーワードがなくなると（S108、NO）、次の分類キーワード候補の比較処理に移行する。未処理のキーワードを選択し、条件キーワード出現回数と主格キーワード出現回数とをカウントする処理（S102～S111）を、分類キーワード候補の全てについて実行し、処理を終了する。

【0086】

上記比較によって、図9に示す集計結果93を得ることができる。

20

【0087】

図9に示す条件キーワード情報91と、主格キーワード情報92との例では、まず、分類キーワード候補「大人だけ」を選択し（S102）、入力された条件キーワードと比較する。条件キーワード「大人だけ」を選択し、比較すると（S105）、完全一致するので、条件キーワード出現回数を、1増やす（S107）。

【0088】

次に、条件キーワード「小さな子供」と比較するが、S106の条件を満たさないので、未処理の条件キーワードが存在するかどうかを調べる（S104）。未処理の条件キーワードが存在しなければ（S104、NO）、主格キーワードとの比較処理を実行する。未処理の主格キーワード「傘」を選択し、分類キーワード候補「大人だけ」と比較する（S109）。この比較の結果、完全一致または部分一致であるというS110の条件を満たさないので、処理S108に戻る。未処理の主格キーワード「大人」が存在するので、「大人」を選択し、分類キーワード候補「大人だけ」と比較する（S109）。

30

【0089】

分類キーワード候補「大人だけ」に、上記比較キーワード「大人」が含まれているので、主格キーワード出現回数を、1増やす。未処理の主格キーワードが無くなるので（S108、NO）、未処理の分類キーワード候補「小さな子供」を選択し（S103）、（S104～S111）の処理を実行する。この結果、分類キーワード候補「小さな子供」について、条件キーワード出現回数が、1増え、1になり、分類キーワード候補「傘」について、主格キーワード出現回数が1増え、1になる。

40

【0090】

また、分類キーワード候補「大人」と、条件キーワード「大人だけ」とを比較すると、部分一致はするが、「大人」に「大人だけ」が含まれていないので、条件キーワード出現回数をカウントすることができない。分類キーワード候補「大人」と、主格キーワード「大人」とを比較すると、完全一致するので、主格キーワード出現回数を、1増やし、1とする。このように処理した結果、図9に示す集計結果93を得ることができる。

【0091】

実施例1では、分類キーワード候補と比較キーワードとが部分一致である場合にも、出現回数をカウントするが、完全一致する場合にのみ、カウントするようにしてもよい。

【0092】

50

このようにして求めた条件キーワード出現回数と主格出現回数とが所定の条件を満たすキーワードを、分類キーワードとして判定する。上記所定の条件は、たとえば、上記式(1)、式(2)に示す条件である。

【0093】

上記式(1)、式(2)において、 $t_h = 0.8$ 、 $t_h = 1$ であるとした場合、各キーワードについての式(1)の値は、以下のようになり、「大人だけ」が、分類キーワードとして抽出される。

【0094】

$$f(\text{大人だけ}) = 0.8 * 1 + 0.2 * 1 = 1.0$$

$$f(\text{小さな子供}) = 0.8 * 1 = 0.8$$

$$f(\text{傘}) = 0.2 * 1 = 0.2$$

$$f(\text{大人}) = 0.2 * 1 = 0.2$$

10

【0095】

実施例1では、説明を簡単にするために、3つの回答文書から、分類キーワードを求めているが、本来、大量の回答文書を入力し、分類キーワードを求める。このために、式(2)の閾値 t_h の値によっては、多くのキーワードが分類キーワードとして判定されることがある。この場合、式(1)の値の上位何件等のように、件数による条件を追加することによって、分類キーワードを判定するようにしてもよい。

【0096】

分類キーワード判定手段15は、分類キーワードと、このキーワードを含む回答文章IDと、文番号の情報とを組みにして、テキスト出力手段40に送る。

20

【0097】

テキスト出力手段40は、分類キーワード判定手段15が出力した分類キーワードを入力し、分類キーワード毎に、回答文書を分類して出力する。分類キーワードとなるキーワード抽出の処理を、文単位に行うので、出力も文単位で行う例について説明する。

【0098】

入力として受け取ったキーワードと、このキーワードを含む回答文章IDと、文番号の情報とに基づいて、質問・回答文書蓄積手段30から、該当する文を検索して取得し、分類キーワード毎に、文の集合を生成する。分類キーワードとこのキーワードを含む文の集合とを一度に出力してもよく、まず、分類キーワードだけを表示し、利用者が選択した分類キーワードを含む文の集合を表示するようにしてもよい。

30

【0099】

つまり、回答文書分類装置100は、質問文書に対する回答文書を、グループを代表するキーワードである分類キーワードを含むグループ毎に分類する回答文書分類装置である。

【0100】

条件パターン蓄積手段11は、単語と品詞との組み合わせによって構成されているパターンであって、所定の条件を示すパターンである条件パターンを蓄積している条件パターン蓄積手段の例である。条件パターン判定手段12は、入力テキスト集合に含まれている文が、上記条件パターン蓄積手段に蓄積されている条件パターンを含んでいるか否かを判定する条件パターン判定手段の例である。条件キーワード抽出手段14は、上記条件パターンを含んでいれば、上記条件パターンの前方に配置されている単語の配列である条件キーワードを抽出する条件キーワード抽出手段の例である。主格キーワード抽出手段14は、上記条件パターンを含んでいなければ、主格になる文節に含まれている単語である主格キーワードを抽出する主格キーワード抽出手段の例である。

40

【0101】

また、分類キーワード判定手段15は、上記条件キーワードおよび上記主格キーワードに含まれるキーワード毎に、着目しているキーワードが上記条件キーワードとして出現する回数である出現回数と、上記主格キーワードとして出現する回数である出現回数とを集計し、上記条件キーワードの出現回数と上記主格キーワードの出現回数とが所定の条件を

50

満たす場合に、上記着目しているキーワードを分類キーワードであると判定する分類キーワード判定手段の例である。

【実施例 2】

【0102】

図 11 は、本発明の実施例 2 である回答文書分類装置 200 を示すブロック図である。

【0103】

回答文書分類装置 200 は、条件キーワード抽出手段 16 と、主格キーワード抽出手段 17 とが抽出した分類キーワード候補を、分類キーワード候補蓄積手段 18 に蓄積し、利用者からの検索要求に応じて、検索結果に含まれている回答文章 ID を持つ文における分類キーワード候補を、分類キーワード蓄積手段 18 から取得し、各分類キーワード候補の条件キーワード出現回数と主格キーワード出現回数とを集計し、分類キーワードを判定する装置である。

10

【0104】

つまり、回答文書分類装置 200 は、テキスト入力手段 21 とテキスト出力手段 40 とに、接続され、条件パターン蓄積手段 11 と、条件パターン判定手段 12 と、条件キーワード抽出手段 16 と、主格キーワード抽出手段 17 と、分類キーワード候補蓄積手段 18 と、分類キーワード判定手段 19 とを有する。なお、実施例 1 における構成要素と同一の構成要素には、同一符号を付してある。

【0105】

テキスト入力手段 21 は、質問・回答文書蓄積手段 30 に蓄積されている未処理の回答文書を定期的に取り得し、条件パターン判定手段 12 に送る。

20

【0106】

条件パターン判定手段 12 と条件パターン蓄積手段 11 とは、実施例 1 の構成と同じであるので、その説明を省略する。

【0107】

条件キーワード抽出手段 16 は、条件パターン判定手段 12 から条件パターンを含む形態素解析結果の文と、条件パターンに一致する箇所とを受け取り、条件パターンに一致する箇所の前方に存在するキーワードを、条件キーワードとして抽出し、記憶装置に記憶する。

【0108】

実施例 1 において、抽出した条件キーワードに、回答文章 ID と文番号とを対応付けて分類キーワード判定手段 15 へ渡したが、実施例 2 では、入力された形態素解析結果に、条件キーワードを示す記号を追加し、分類キーワード候補蓄積手段 18 に記録する。なお、上記のように、入力された形態素解析結果に、条件キーワードを示す記号を追加する理由は、条件キーワードを明示するためである。

30

【0109】

図 12 は、実施例 2 における分類キーワード蓄積手段 18 の蓄積例を示す図である。

【0110】

条件キーワードを示す記号 122、123 は、分類キーワード候補蓄積手段 18 に蓄積されている分類キーワード候補が条件キーワードであることを示す記号であり、分類キーワード候補の行末に、記号「X」が付与されている。記号「X」が複数行に渡って付与されていれば、単語ではなく文節が抽出されていることを示す。

40

【0111】

主格キーワード抽出手段 17 は、条件パターン判定手段 12 から、条件パターンを含まないと判定された文を受け取り、係り受け解析を行い、この解析結果から、文の中に主格となる文節が存在するかどうかを調べ、存在すれば、主格となる文節に含まれているキーワード(単語、名詞、名詞句)を、主格キーワードとして抽出し、記憶装置に記憶する。

【0112】

実施例 1 において、抽出した主格キーワードに、回答文章 ID と文番号とを対応付け、分類キーワード判定手段 15 へ渡したが、実施例 2 では、入力された形態素解析結果に、

50

主格キーワードを示す記号を追加し、分類キーワード候補蓄積手段 18 に記録する。主格キーワードを示す記号 121、124 は、分類キーワード候補蓄積手段 18 に蓄積されている分類キーワードが主格キーワードであることを示す記号であり、分類キーワード候補の行末に、記号 Y が付与されている。

【0113】

テキスト入力手段 21 において、処理を定期的に行う度に、条件パターン判定手段 12、条件キーワード抽出手段 16、主格キーワード抽出手段 17 が処理し、分類キーワード候補蓄積手段 18 に、分類キーワード候補を含む文を蓄積する。

【0114】

次に、利用者が、キーワード検索手段 50 に検索キーワードを入力すると、キーワード検索手段 50 は、質問・回答文書蓄積手段 30 を検索し、検索キーワードを含む質問文書について回答している回答文書の集合を取得する。取得した回答文書集合の回答文章 ID の集合を、分類キーワード判定手段 19 に渡す。

10

【0115】

分類キーワード判定手段 19 は、処理対象となるキーワード検索手段 50 から、回答文章 ID を受け取ると、回答文章 ID が一致する文を、分類キーワード候補蓄積手段 18 から取得する。各文から、条件キーワードまたは主格キーワードを取得し、「条件キーワード 回答文章 ID 文番号」の組からなる条件キーワード情報 91 の集合と、「主格キーワード 回答文章 ID 文番号」の組からなる主格キーワード情報 92 の集合とを取得し、図 10 に示す分類キーワード判定処理を実行し、分類キーワードを判定する。判定した結果、分類キーワードと、このキーワードを含む回答文章 ID と、文番号の情報とを組みにして、テキスト出力手段 40 に送る。

20

【0116】

実施例 1 では、テキスト入力手段 20 から、回答文章 ID が与えられてから、分類キーワードを出力するまでに、形態素解析処理、係り受け処理、分類キーワード判定処理の時間を加算した処理時間が必要である。これらの処理は、文章数に比例する処理であるので、従来技術で説明したクラスタリング処理（文章数の 2 乗に比例する）に比較すれば高速である。

【0117】

実施例 2 では、分類キーワード候補蓄積手段 18 が、分類キーワード候補を予め抽出して蓄積するので、分類キーワード判定処理のみの時間で、分類キーワードを出力することができ、大変高速である。

30

【0118】

分類キーワード候補蓄積手段 18 は、回答文書の文章毎に、条件キーワードであることを示す情報または主格キーワードであることを示す情報が追加されている解析結果と、回答文書を一意に示す回答文章 ID とを対応付けた分類キーワード候補を蓄積する分類キーワード候補蓄積手段の例である。

【0119】

分類キーワード判定手段 19 は、回答文章 ID が複数与えられた際に、回答文章 ID に該当する回答文書の解析結果の集合を、上記分類キーワード候補蓄積手段から取得し、上記解析結果の集合に含まれている条件キーワードおよび主格キーワードに含まれるキーワード毎に、着目しているキーワードが上記条件キーワードとして出現する回数である出現回数と、上記主格キーワードとして出現する回数である出現回数とを集計し、上記条件キーワードの出現回数と上記主格キーワードの出現回数とが所定の条件を満たす場合に、上記着目しているキーワードを分類キーワードであると判定する分類キーワード判定手段の例である。

40

【0120】

また、上記実施例において、上記条件キーワード抽出手段は、上記条件パターンを含む文節から、条件パターンを除いた語句に、条件パターンを含む文節へ係る 0 個以上の文節を追加した語句を、条件キーワードとして抽出する手段である。また、上記主格キーワー

50

ド抽出手段は、主格になる文節に含まれている名詞句に、主格になる文節へ係る 0 個以上の文節を追加した語句を、上記条件キーワードとして抽出する手段である。そして、上記分類キーワード判定手段は、上記条件キーワードとしての出現回数が多いほど、上記分類キーワードとして判定され易くなる条件を用いて、上記分類キーワードを判定する手段である。

【 0 1 2 1 】

上記実施例において、手段を工程に置き換えれば、上記実施例を方法の発明として把握することができる。つまり、上記実施例は、質問文書に対する回答文書を、グループを代表するキーワードである分類キーワードを含むグループ毎に分類する回答文書分類方法において、単語と品詞との組み合わせによって構成されているパターンであって、所定の条件を示すパターンである条件パターンが蓄積されている条件パターン蓄積手段に蓄積されている条件パターンを、入力テキスト集合に含まれている文が含んでいるか否かを、条件パターン判定手段が判定し、記憶装置に記憶する条件パターン判定工程と、上記条件パターンを含んでいれば、上記条件パターンの前方に配置されている単語の配列である条件キーワードを、条件キーワード抽出手段が抽出し、記憶装置に記憶する条件キーワード抽出工程と、上記条件パターンを含んでいなければ、主格になる文節に含まれている単語である主格キーワードを、主格キーワード抽出手段が抽出し、記憶装置に記憶する主格キーワード抽出工程と、上記条件キーワードおよび上記主格キーワードに含まれるキーワード毎に、着目しているキーワードが上記条件キーワードとして出現する回数である出現回数と、上記主格キーワードとして出現する回数である出現回数とを集計し、上記条件キーワードの出現回数と上記主格キーワードの出現回数とが所定の条件を満たす場合に、上記着目しているキーワードを分類キーワードであると、分類キーワード判定手段が判定し、記憶装置に記憶する分類キーワード判定工程とを有する回答文書分類方法の例である。

【 0 1 2 2 】

また、上記実施例は、質問文書に対する回答文書を、グループを代表するキーワードである分類キーワードを含むグループ毎に分類する回答文書分類方法において、単語と品詞との組み合わせによって構成されているパターンであって、所定の条件を示すパターンである条件パターンを蓄積している条件パターン蓄積手段に蓄積されている条件パターンを、回答文書集合に含まれている文が含んでいるか否かを、条件パターン判定手段が判定し、記憶装置に記憶する条件パターン判定工程と、上記条件パターンを含んでいれば、上記条件パターンの前方に配置されている単語の配列である条件キーワードを、条件キーワード抽出手段が抽出し、記憶装置に記憶する条件キーワード抽出工程と、上記条件パターンを含んでいなければ、主格になる文節に含まれている単語である主格キーワードを、主格キーワード抽出手段が抽出し、記憶装置に記憶する主格キーワード抽出工程と、回答文書の文章毎に、条件キーワードであることを示す情報または主格キーワードであることを示す情報が追加されている解析結果と、回答文書を一意に示す回答文章 ID とを対応付けた分類キーワード候補とを、分類キーワード候補蓄積手段が蓄積し、記憶装置に記憶する分類キーワード候補蓄積工程と、回答文章 ID が複数与えられた際に、回答文章 ID に該当する回答文書の上記解析結果の集合を、上記分類キーワード候補蓄積工程から取得し、上記解析結果の集合に含まれている条件キーワード毎にまたは主格キーワード毎に、着目しているキーワードが上記条件キーワードとして出現する回数である出現回数と、上記主格キーワードとして出現する回数である出現回数とを集計し、上記条件キーワードの出現回数と上記主格キーワードの出現回数とが所定の条件を満たす場合に、上記着目しているキーワードを分類キーワードであると、分類キーワード判定が判定し、記憶装置に記憶する分類キーワード判定工程とを有する回答文書分類方法の例である。

【 0 1 2 3 】

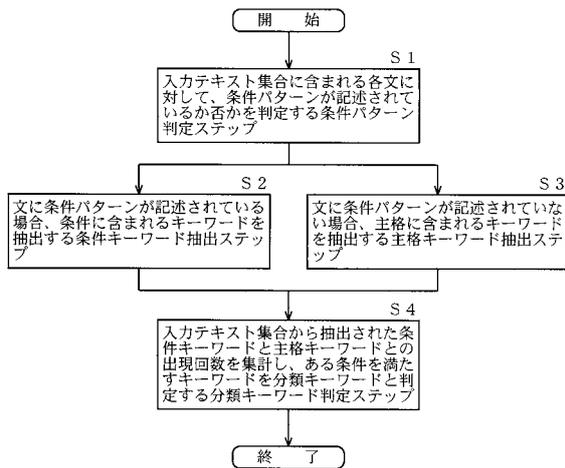
さらに、上記実施例をプログラムの発明として把握することができる。つまり、上記実施例は、請求項 1 ~ 請求項 5 のいずれか 1 項に記載の回答文書分類装置を構成する各手段としてコンピュータを機能させるプログラムの例である。

【 符号の説明 】

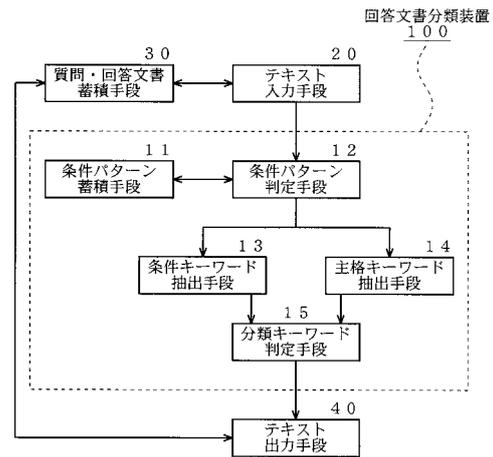
【 0 1 2 4 】

- 1 0 0 ... 回答文書分類装置、
- 1 1 ... 条件パターン蓄積手段、
- 1 2 ... 条件パターン判定手段、
- 1 3 ... 条件キーワード抽出手段、
- 1 4 ... 主格キーワード抽出手段、
- 1 5 ... 分類キーワード判定手段、
- 2 0 ... テキスト入力手段、
- 3 0 ... 質問・回答文書蓄積手段、
- 4 0 ... テキスト出力手段、
- 2 0 0 ... 回答文書分類装置、
- 1 6 ... 条件キーワード抽出手段、
- 1 7 ... 主格キーワード抽出手段、
- 1 8 ... 分類キーワード候補蓄積手段、
- 1 9 ... 分類キーワード判定手段、
- 5 0 ... キーワード検索手段。

【 図 1 】



【 図 2 】



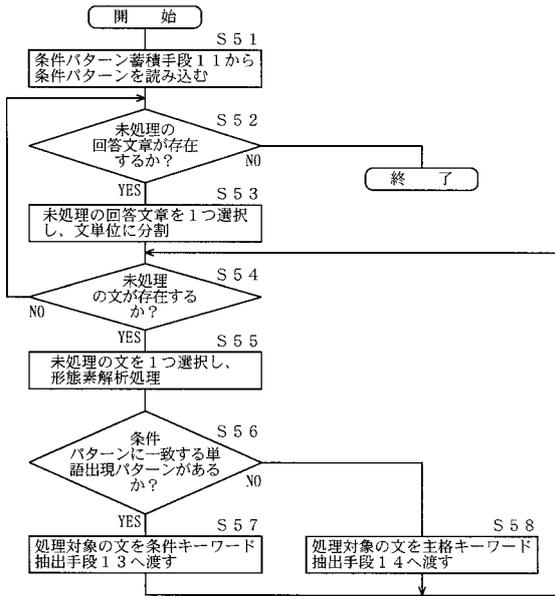
【 図 3 】

条件パターン蓄積手段 11 が蓄積している条件パターンの例

接続接尾辞：ナラ
 判定詞：ナラ
 判定詞：ナラバ
 判定詞：デ、動詞語幹：ア、動詞活用語尾：レ、動詞接尾辞：バ
 格助詞：ノ、名詞：バアイ
 判定詞：デス、引用助詞：ト

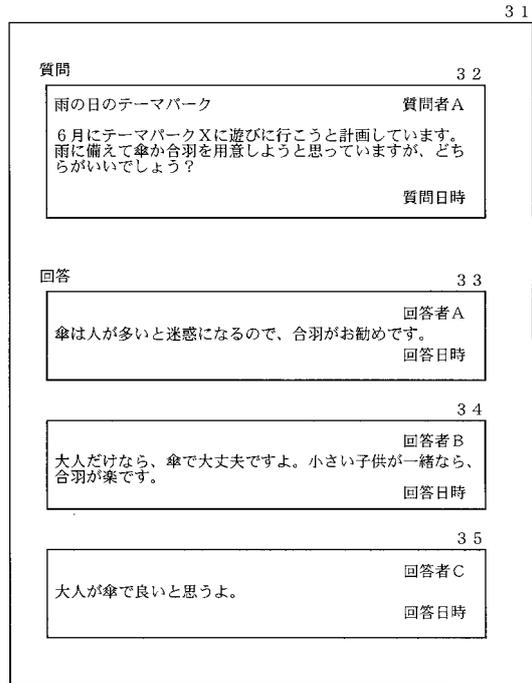
【図4】

条件パターン判定手段1.2の動作



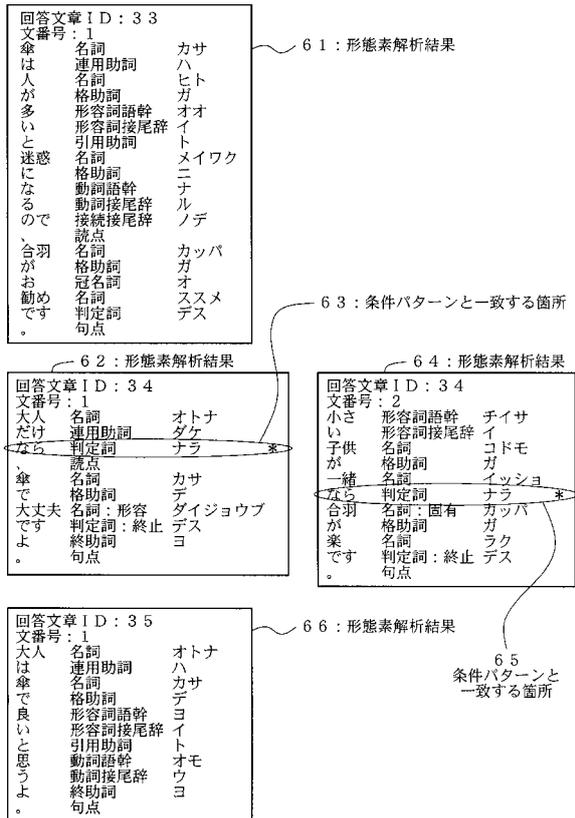
【図5】

Q&Aコミュニティサイトの構成



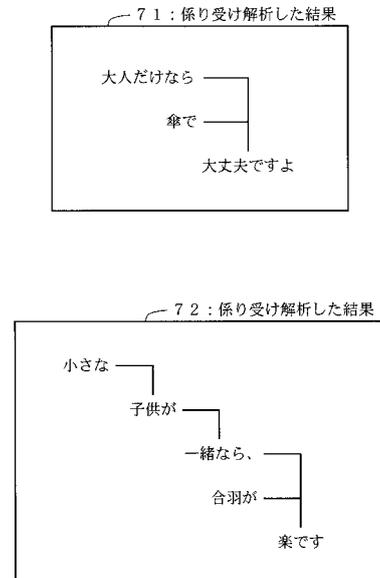
【図6】

条件パターン判定手段1.2が判定した形態素解析結果の例



【図7】

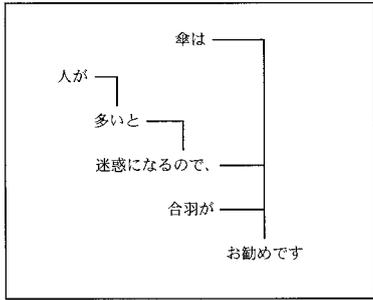
条件キーワード抽出手段1.3の動作



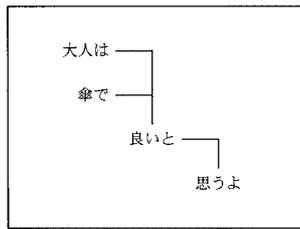
【図 8】

主格キーワード抽出手段 14 の動作

(1)



(2)



【図 9】

分類キーワード判定手段 15 の動作説明

91: 条件キーワード情報

大人だけ	3	4	1
小さな子供	3	4	2

92: 主格キーワード情報

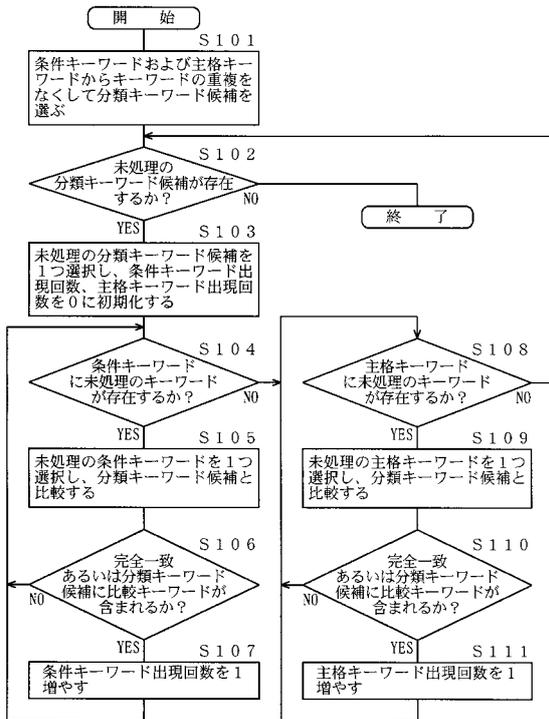
傘	3	3	1
大人	3	5	1

集計結果
93

分類キーワード候補	条件キーワード出現回数	主格キーワード出現回数	条件キーワード	主格キーワード
大人だけ	1	1	大人だけ 3 4 1	大人 3 5 1
小さな子供	1	0	小さな子供 3 4 2	—
傘	0	1	—	傘 3 3 1
大人	0	1	—	大人 3 5 1

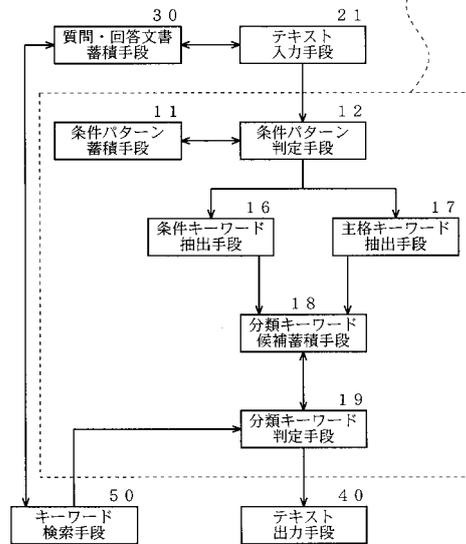
【図 10】

分類キーワード判定手段 15 の動作



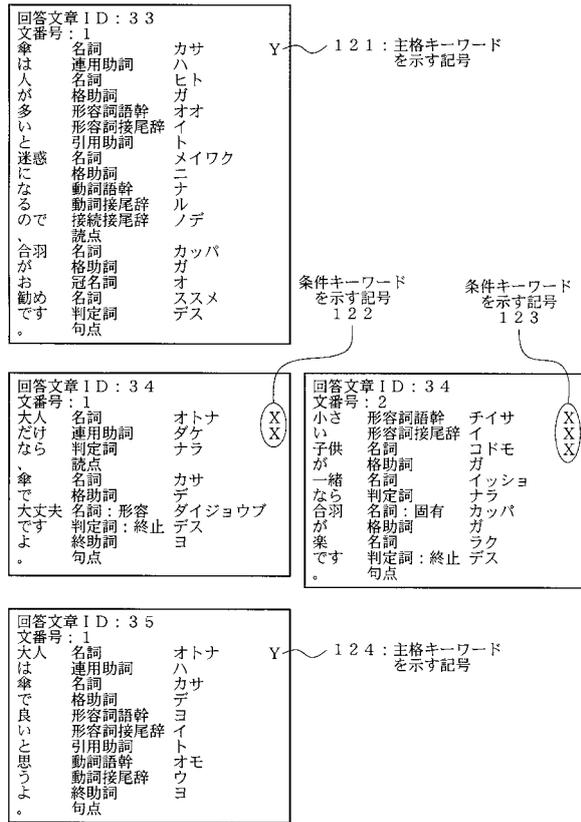
【図 11】

回答文書分類装置
200



【図 12】

分類キーワード候補蓄積手段18における蓄積例



フロントページの続き

- (56)参考文献 特開平05-120345(JP,A)
特開平01-121928(JP,A)
特開平11-203318(JP,A)
特開2000-259666(JP,A)
別所 克人 他, 自然言語検索システムにおける分野推論方式, 電子情報通信学会論文誌, 日本,
電子情報通信学会, 1998年 6月25日, Vol.J81-D-II, No.6, PP.1317-1327.
佐々木 裕 他, SVMを用いた学習型質問応答システムS A I Q A - I I, 情報処理学会論文誌,
日本, 情報処理学会, 2004年 2月15日, Vol.45, No.2, PP.635-646.

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