

(19) 日本国特許庁(JP)

(12) 特許公報(B2)

(11) 特許番号

特許第5167546号  
(P5167546)

(45) 発行日 平成25年3月21日(2013.3.21)

(24) 登録日 平成25年1月11日(2013.1.11)

(51) Int.Cl. F I  
**G06F 17/30 (2006.01)** G O 6 F 17/30 3 5 0 C  
 G O 6 F 17/30 1 7 0 A

請求項の数 21 (全 66 頁)

(21) 出願番号	特願2008-530812 (P2008-530812)	(73) 特許権者	504132272
(86) (22) 出願日	平成19年3月16日 (2007.3.16)		国立大学法人京都大学
(86) 国際出願番号	PCT/JP2007/055448		京都府京都市左京区吉田本町36番地1
(87) 国際公開番号	W02008/023470	(74) 代理人	100078868
(87) 国際公開日	平成20年2月28日 (2008.2.28)		弁理士 河野 登夫
審査請求日	平成22年3月9日 (2010.3.9)	(74) 代理人	100114557
(31) 優先権主張番号	特願2006-224563 (P2006-224563)		弁理士 河野 英仁
(32) 優先日	平成18年8月21日 (2006.8.21)	(72) 発明者	白松 俊
(33) 優先権主張国	日本国 (JP)		京都府京都市左京区吉田本町 京都大学大学院情報学研究科内
		(72) 発明者	駒谷 和範
			京都府京都市左京区吉田本町 京都大学大学院情報学研究科内

最終頁に続く

(54) 【発明の名称】 文単位検索方法、文単位検索装置、コンピュータプログラム、記録媒体及び文書記憶装置

(57) 【特許請求の範囲】

【請求項1】

自然言語からなる複数の文書データが記憶されている文書集合を用い、該文書集合から取得した文書データを一又は複数の文からなる文単位に分別しておく一方、言葉を順次受け付け、受け付けた言葉に基づいて前記文書集合から分別してある文単位を検索する文単位検索方法において、

文書データ中に連なる文単位夫々に、該文単位及び先行文脈に基づき求められる該文単位での顕現性を表わす重み値が付与された複数の単語からなる重み付き単語群を対応付けて予め記憶しておくステップと、

言葉を受け付ける都度、該言葉に、該言葉及び先行文脈に基づき求められる該言葉での顕現性を表わす重み値が付与された複数の単語からなる重み付き単語群を対応付けるステップと、

受け付けた言葉に対応付けた重み付き単語群と類似する重み付き単語群が対応付けて記録されている文単位を、前記文書集合から抽出する類似文単位抽出ステップと、

抽出した文単位を出力するステップと  
 を含むことを特徴とする文単位検索方法。

【請求項2】

前記類似文単位抽出ステップは、

受け付けた言葉に対応付けた重み付き単語群の内の複数の単語の重み値の分布と、予め分別された文単位に対応付けられている重み付き単語群の内の複数の単語の重み値の分布

とが、所定の条件を満たすか否かを判断するステップと、

所定の条件を満たすと判断された重み付き単語群が対応付けられている文単位を抽出するステップと

を含むことを特徴とする請求項 1 に記載の文単位検索方法。

【請求項 3】

前記類似文単位抽出ステップは、

予め分別された文単位から、受け付けた言葉に対応付けた重み付き単語群と同一の単語を含む単語群が対応付けられた文単位を抽出するステップと、

受け付けた言葉と抽出した文単位とで、対応付けられた単語群の内の同一の単語毎に重み値の差分を算出するステップと、

抽出した文単位に、算出した差分が小さい順に優先順位を付与するステップと

を含み、

抽出した文単位を、優先順位に基づいて出力する

ことを特徴とする請求項 1 又は 2 に記載の文単位検索方法。

【請求項 4】

前記重み付き単語群を、各単語を 1 次元とし、単語毎に付与される重み値の大きさを各単語に対応する次元方向の要素として持つ多次元ベクトルとして算出するステップを含み、

前記類似文単位抽出ステップは、

分別した文単位毎に記憶してある前記多次元ベクトルと、受け付けた言葉に対応付けた前記多次元ベクトルとの距離を算出するステップと、

文単位に、算出した距離が短い順に優先順位を付与するステップと

を含み、

付与された優先順位に従って出力する

ことを特徴とする請求項 1 又は 2 に記載の文単位検索方法。

【請求項 5】

文単位又は受け付けた言葉に重み付き単語群を対応付ける際、

各単語が、前記文単位又は前記言葉よりも後続の文単位又は言葉に出現する又は参照される参照確率を算出する参照確率算出ステップを含み、

算出した参照確率を各単語の重み値として付与する

ことを特徴とする請求項 1 乃至 4 のいずれかに記載の文単位検索方法。

【請求項 6】

前記参照確率算出ステップは、

前記各単語が先行の文単位を含む複数の文単位に出現するパターン、又は前記単語を先行の文単位から参照するパターンを含む特徴パターンを特定するステップと、

前記文書集合から取得された文書データ中で、前記特徴パターンと同一の特徴パターンが特定される単語が、後続の文単位で出現する又は参照される割合を算出するステップとを含み、

算出した割合を参照確率とする

ことを特徴とする請求項 5 に記載の文単位検索方法。

【請求項 7】

前記文書集合から抽出される単語毎に、該単語が先行の文単位を含む複数の文単位に出現するパターン、又は前記単語を先行の文単位から参照するパターンを含む特徴パターンを特定する特定ステップと、

特定した特徴パターンと同一の特徴パターンが特定される単語が、前記文書データ中で後続の文単位で出現したか又は参照されたかを判定する判定ステップと、

特定した特徴パターンと、該特徴パターンで特定される単語に対して判定した結果との回帰分析を行って前記参照確率に対する前記特徴パターンの回帰係数を算出する回帰ステップと

を含み、

10

20

30

40

50

文単位に重み付き単語群を対応付けて記憶しておく際、又は受け付けた言葉に重み付き単語群を対応付ける際、

前記参照確率算出ステップは、

前記文単位又は言葉毎に、該文単位又は言葉での単語の特徴パターンを特定し、

特定した特徴パターンに対する前記回帰係数を使用して参照確率を算出する

ことを特徴とする請求項 5 に記載の文単位検索方法。

【請求項 8】

文単位に対しては、書き言葉からなる第 1 文書集合から取得された文書データ中で前記割合を算出し、

受け付けた言葉に対しては、話し言葉からなる第 2 文書集合から取得された文書データ中で前記割合を算出する

ことを特徴とする請求項 6 に記載の文単位検索方法。

10

【請求項 9】

書き言葉からなる第 1 文書集合及び話し言葉からなる第 2 文書集合夫々について、

前記特定ステップ、前記判定ステップ及び前記回帰ステップを実行しておき、

前記参照確率算出ステップは、

前記文単位で特定した単語の特徴パターンに対しては、第 1 文書集合について実行した前記回帰ステップにより算出された回帰係数を使用して参照確率を算出し、

前記受け付けた言葉で特定した単語の特徴パターンに対しては、第 2 文書集合について実行した前記回帰ステップで算出された回帰係数を使用して参照確率を算出する

ことを特徴とする請求項 7 に記載の文単位検索方法。

20

【請求項 10】

前記特徴パターンは、

前記単語を先行の文単位又は言葉から参照している場合の前記先行の文単位又は言葉から前記単語が含まれる文単位又は言葉までの、文単位又は言葉の数、

前記単語が出現又は参照されている直近の先行の文単位又は言葉における前記単語の係り受け情報、

前記単語が含まれる文単位又は言葉までに出現した又は参照された回数、

前記単語が出現又は参照されている直近の先行の文単位又は言葉における前記単語の名詞区別、

30

前記単語が出現又は参照されている直近の先行の文単位又は言葉中で前記単語が主題であるか否か、

前記単語が出現又は参照されている直近の先行の文単位又は言葉中で前記単語が主語であるか否か、

前記単語が含まれる文単位又は言葉における人称、

及び、

前記単語が含まれる文単位又は言葉における品詞情報、

の内の一又は複数を含む情報で特定される

ことを特徴とする請求項 6 乃至 9 のいずれかに記載の文単位検索方法。

【請求項 11】

40

前記特徴パターンは、

前記単語を先行の文単位又は言葉から参照している場合の前記先行の文単位又は言葉から前記単語が含まれる文単位又は言葉までに対応する時間、

前記単語が出現又は参照されている直近の先行の文単位又は言葉中で前記単語に対応する発話速度、

及び、

前記単語が出現又は参照されている直近の先行の文単位又は言葉中で前記単語に対応する音声の周波数

の内の一又は複数を含む情報で特定される

ことを特徴とする請求項 6 乃至 10 のいずれかに記載の文単位検索方法。

50

## 【請求項 1 2】

前記文章集合から抽出される単語の内の一の単語について、

前記分別された文単位に対応付けられている重み付き単語群の内から、前記一の単語が含まれる単語群であり、且つ前記一の単語の重み値が所定値以上である単語群を抽出する第 1 ステップと、

該第 1 ステップで抽出した単語群の各単語の重み値を単語毎に統合した値を、前記一の単語の各単語への関連度として付与した関連単語群を作成する第 2 ステップと、

作成した関連単語群を前記一の単語に対応付けて記憶する第 3 ステップと、

前記抽出された単語夫々について前記第 1 ステップ乃至第 3 ステップを予め実行するステップと、

文単位毎又は受け付けた言葉毎に対応付けられた重み付き単語群の各単語の重み値夫々を、各単語に対応付けて記憶されている前記関連単語群の各単語の関連度を使用して付与し直す関連度付加ステップと

を含むことを特徴とする請求項 1 乃至 1 1 のいずれかに記載の文単位検索方法。

10

## 【請求項 1 3】

前記第 2 ステップは、

前記抽出した単語群について、各単語群に含まれる各単語の重み値に、前記一の単語の重み値で重み付けした総和を算出するステップと、

算出した総和を平均化するステップと、

作成する関連単語群の各単語の前記関連度として、各単語の重み値の平均化された総和を付与するステップと

を含むことを特徴とする請求項 1 2 に記載の文単位検索方法。

20

## 【請求項 1 4】

前記関連度付加ステップは、

文単位毎又は受け付けた言葉毎に対応付けられた重み付き単語群の各単語について、

各単語に対応付けて記憶されている前記関連単語群に含まれる各単語の関連度を、前記重み付き単語群の各単語の重み値に乗算するステップと、

乗算結果に基づいて前記重み付き単語群の各単語の重み値として付与し直すステップと

を含むことを特徴とする請求項 1 2 又は 1 3 に記載の文単位検索方法。

30

## 【請求項 1 5】

各単語夫々についての前記関連単語群を、各単語を 1 次元とし、単語毎に付与される関連度の大きさを各単語に対応する次元方向の要素として持つ多次元の関連度ベクトルとして算出するステップと

を含み、

前記関連度付加ステップは、

分別した文単位毎に記憶してある前記多次元ベクトルを、各単語の関連度ベクトルの列によって変換する

ことを特徴とする請求項 1 2 乃至 1 4 のいずれかに記載の文単位検索方法。

## 【請求項 1 6】

自然言語からなる複数の文書データが記憶されている文書集合を用い、言葉を受け付け、受け付けた言葉に基づいて前記文書集合を検索する文単位検索方法において、

前記文書集合から得られる文書データを一又は複数の文からなる文単位に分別しておくステップ、

分別した文単位毎に、該文単位に出現する単語、又は、文書データ中の先行の文単位から参照する単語を抽出するステップ、

前記文単位に対して抽出した単語毎に、各文単位における特徴を特定して記憶しておくステップ、

分別した文単位毎に、該文単位に対して抽出した単語が該文単位及び先行の文単位で出現する場合の前記特徴の組み合わせのパターン、又は先行の文単位から参照する場合の参照のパターンを含む特徴パターンを特定するステップ、

40

50

特定した特徴パターンと、該特徴パターンで特定された単語が後続の文単位で出現又は参照されたか否かとを記憶しておくステップ、

前記文書集合から得られる文書中の文単位全体に対し、一の特徴パターンで特定される単語が後続の文単位で出現又は参照される参照確率の回帰分析を行って特徴パターンに対応する回帰係数を得る回帰学習を実行するステップ、

分別した文単位毎に、

文書データ中で先行の文単位から各文単位に至るまでに抽出された各単語について、前記文単位で特定される特徴パターンに対応する前記回帰係数を使用し、前記単語の前記参照確率を算出するステップ、

算出した参照確率を夫々付与した重み付き単語群を対応付けて予め記憶しておくステップ、

言葉を受け付けた場合、受け付けた順に言葉を記憶するステップ、

言葉を受け付けた場合、

受け付けた言葉に出現する単語又は前記言葉よりも先に受け付けた言葉から参照する単語を抽出するステップ、

抽出した各単語の前記受け付けた言葉における特徴を特定するステップ、

先に受け付けた言葉で出現する場合の特徴の組み合わせのパターン、又は先に受け付けた言葉から参照する場合の参照のパターンを含む特徴パターンを特定するステップ、

特定された特徴パターンに対応する前記回帰係数を使用して、前記単語の前記参照確率を算出するステップ、

算出した参照確率を夫々付与した重み付き単語群を前記言葉に対応付けるステップ、

前記受け付けた言葉と、予め分別されてある文単位とで、対応付けられている重み付き単語群の内の同一の単語毎に付与されている参照確率の差分を算出するステップ、

予め分別されてある文単位に、前記参照確率の差分が小さい順に優先順位を付与するステップ、及び、

前記文単位を付与された優先順位に基づいて出力するステップ

を含むことを特徴とする文単位検索方法。

#### 【請求項 17】

自然言語からなる複数の文書データが記憶されている文書集合から文書データを取得する手段と、言葉を順次受け付ける手段とを備え、受け付けた言葉に基づいて前記文書集合を検索する文単位検索装置において、

取得した文書データを一又は複数の文からなる文単位に分別する手段と、

取得した文書データ中に連なる文単位夫々に、該文単位及び先行文脈に基づき求められる該文単位での顕現性を表わす重み値が付与された複数の単語からなる重み付き単語群を対応付けて記憶する手段と、

言葉を受け付けた場合に受け付けた順に記憶する手段と、

新たに言葉を受け付ける都度、該言葉に、該言葉及び該先行文脈に基づき求められる該言葉での顕現性を表わす重み値が付与された複数の単語からなる重み付き単語群を対応付ける手段と、

予め分別された文単位から、受け付けた言葉に対応付けた重み付き単語群と類似する重み付き単語群が対応付けて記録されている文単位を抽出する手段と、

抽出した文単位を出力する手段と

を備えることを特徴とする文単位検索装置。

#### 【請求項 18】

自然言語からなる複数の文書データが記憶されている文書集合から、文書データを取得することが可能であるコンピュータを、言葉を順次受け付ける手段と、受け付けた言葉に基づいて前記文書集合を検索する手段として機能させることができるコンピュータプログラムにおいて、

取得した文書データを一又は複数の文からなる文単位に分別する手段、

取得した文書データ中に連なる文単位夫々に、該文単位及び先行文脈に基づき求められ

10

20

30

40

50

る該文単位での顕現性を表わす重み値が付与された複数の単語からなる重み付き単語群を対応付けて記憶する手段、

言葉を受け付けた場合に受け付けた順に記憶する手段、

新たに言葉を受け付ける都度、該言葉に、該言葉及び先行文脈に基づき求められる該言葉での顕現性を表わす重み値が付与された複数の単語からなる重み付き単語群を対応付ける手段、及び、

予め分別された文単位から、受け付けた言葉に対応付けた重み付き単語群と類似する重み付き単語群が対応付けて記録されている文単位を抽出する手段

として機能させることを特徴とするコンピュータプログラム。

【請求項 19】

請求項 18 に記載のコンピュータプログラムを記録した、コンピュータで読み取り可能な記録媒体。

【請求項 20】

自然言語からなる複数の文書データを記憶する手段と、記憶した文書データを、文書データの先頭から順に一又は複数の文からなる文単位に分別する手段とを備え、分別した文単位毎に、該文単位に出現する単語又は先行する文単位から参照する単語が抽出してあり、分別した文単位毎に抽出した単語が記憶してある文書記憶装置において、

文書データ中に連なる文単位毎に、複数の単語が、該文単位よりも後続の文単位に出現するか又は参照される参照確率を算出する手段と、

前記文単位夫々に、該文単位での顕現性を表わす重み値として前記参照確率が付与された前記複数の単語からなる重み付き単語群を対応付けて記憶する手段と

を備えることを特徴とする文書記憶装置。

【請求項 21】

抽出されてある単語の内の一の単語について、

文単位夫々に対応付けられている重み付き単語群の内から、前記一の単語が含まれる単語群であり、且つ前記一の単語の重み値が所定値以上である単語群を抽出する抽出手段と

、  
該抽出手段が抽出した単語群の各単語の重み値を単語毎に統合した値を、前記一の単語の各単語への関連度として付与した関連単語群を作成する作成手段と、

作成した関連単語群を前記一の単語に対応付けて記憶する記憶手段と

を備え、

前記抽出されてある単語夫々について前記抽出手段、前記作成手段及び前記記憶手段の処理を実行するようにしてあり、各単語に対応付けて夫々の関連単語群を記憶するようにしてあること

を特徴とする請求項 20 に記載の文書記憶装置。

【発明の詳細な説明】

【技術分野】

【0001】

本発明は、検索のためにユーザから受け付けたテキスト、音声等の言葉に基づいて、多数の文書データ記憶されている文書集合からの検索を行う検索方法に関する。特に、文脈の流れの中で意味が動的に変化する文書中の意味のまとまりの単位である文単位から、受け付けた言葉と意味合いが類似する文単位を直接的に検索することができる文単位検索方法、文単位検索装置、コンピュータを前記文単位検索装置として機能させるコンピュータプログラム、該コンピュータプログラムを記録したコンピュータ読み取り可能な記録媒体、及び文書記憶装置に関する。

【背景技術】

【0002】

インターネット上で提供される各種サービスには、ユーザによって入力されたキーワード又は文に基づいて、インターネットで公開されている文書から関連する文書を検索し、一覧にして出力する文書検索サービスがある。

10

20

30

40

50

## 【 0 0 0 3 】

従来の文書検索サービスには、以下のようなものがある。インターネットで公開されている文書を自動的に集めて記憶し、夫々の文書毎に、文書中に出現する単語を文書中での出現確率と共に記憶しておき、キーワード又は文等の言葉を受け付けた場合に、記憶した文書集合から受け付けたキーワード又は文に含まれる単語の出現確率の高い順に優先順位を付与して文書を抽出し、抽出した文書から、当該単語が含まれる文又は段落を出力する。

## 【 0 0 0 4 】

文書検索サービスを利用するユーザは、知りたい情報を検索するために関連するキーワードを自分で考える必要がある。最近の文書検索サービスでは、自然文を入力文として受け付け、入力文を形態素解析し、入力文のキーワードを識別して検索要求を自動的に作成することができる場合もある。

10

## 【 0 0 0 5 】

また、文書検索サービスでは通常、自然文の入力を受け付ける場合でも、入力文に含まれる単語を抽出し、抽出した単語が含まれている文書を検索結果として出力する。したがって、ユーザは、目的の検索結果を得るために入力するキーワードに関連するキーワード又は入力するキーワードの意味付けが変化する単語を更に入力して絞込みをさせる必要があった。例えば、単に「大統領」では、どの国の大統領なのかは不明であるため、「大統領、アメリカ」とキーワードを付加する必要がある。更にアメリカの大統領の何を調べたいかによって、「大統領、アメリカ、出身」、「大統領、アメリカ、政策」等、検索結果を得やすくするための情報を考える必要がある。

20

## 【 0 0 0 6 】

したがって、ユーザが得たいと考える検索結果を実際に得るためには、ユーザはキーワードの組み合わせを考え、何回か試行することが必要になる。例えば、ユーザが「アメリカの大統領は、他の国との間で経済面の問題が発生した場合どのような対策をとるのか」という情報を知りたい場合であっても、「アメリカ、大統領、経済」では検索結果が大量に出力され、大量に出力された検索結果からユーザは文書を選択しなければならない。そこで例えば、「政策」というキーワードを付加して絞込み、「アメリカ、大統領、経済、政策」というキーワードを入力する。この場合、「政策」という言葉が意味の広い上位概念であっても、「政策」というキーワード自体で絞込みをすることになるため、内容としては経済政策についての論述が記載された文書も、「政策」という言葉の出現頻度が低い文書は漏れてしまうことがある。このように、ユーザが検索の目的を達するためのキーワードを考えて試行することで検索結果を得るのは難しい。付加的な情報を入力する度に、本来の検索の目的から、検索結果の内容が離れていく場合もある。

30

## 【 0 0 0 7 】

また、上述の例でユーザが知りたいのは、経済面での政策であって、しかも国際的な政策についてである。ユーザの入力が自然文によるものであっても、「アメリカ、大統領、他の国、経済、問題、発生、場合、対策」の単語の何れの単語が一番重要であるのかは、人間が読む場合は把握できるが、装置又はコンピュータが扱う情報量として定量的に表現することは難しい。したがって、キーワードは全て含んでいるものの、「アメリカの経済の問題と他国の大統領の対策」とについて論述された文書が出力されることも想定できる。

40

## 【 0 0 0 8 】

さらに、検索対象である文書が非常に長い場合は、その文書の中で文脈が動的に変化しているにも拘わらず、その文書を一単位として出現する単語に基づいた検索がされる。したがって、アメリカの大統領の歴史と、他の国の大統領の歴史と、各国の経済のしくみと、各国での失業対策についての内容とが章に分けられて記載されている文書が存在する場合、検索のキーワードをほとんど含むために検索結果として出力される。実際にはそれらの章が文脈的に繋がっていない場合でも、キーワードを含む文又は段落を夫々部分的に抽出した結果が出力されてしまう。そのため、その抽出された部分に至るまでの先行文脈の

50

影響を含む意味と、ユーザの意識の上での検索意図とが、意味的にマッチするか否かは量り得ない。

【0009】

一方、検索対象である文書に、検索のために入力したキーワードは頻繁に出現してはいないにも拘わらず、入力したキーワードが文脈上重要な意味を持って含まれている場合がある。例えば、主題となる単語ほど指示代名詞又はゼロ代名詞で表現される。したがって、知りたい情報を検索するユーザは、検索のために入力したキーワードが指示代名詞又はゼロ代名詞で表現されている文又は段落こそ、検索結果として得たい情報である場合が考えられる。しかしながら、実際の出現頻度で検索結果に優先順位を付与する場合、ユーザが入力したキーワードの出現頻度が低いために絞込みによって候補から除かれ、検索結果として出力されない。

10

【0010】

そこで、文書中の単語を抽出し、当該単語の品詞情報、単語間の係り受け情報、更に指示代名詞又はゼロ代名詞と照応関係にある単語を明示した情報を、文書を形態素解析等により解析した結果に付加して記憶させておき、記憶させた情報に基づいて装置又はコンピュータによる文書の検索、質問応答、機械翻訳を実現する技術が提案されている（非特許文献1）。

【0011】

単語間の係り受け又は照応等の関係は、自然文であるがために文節の順序が複雑であり、人間が読む場合は意味を判別できても機械的に認識することが難しい。そこで、非特許文献1に記載されている技術では、単語間の係り受け又は照応等の関係をタグによって文又は句毎の情報として文書データに付加して記憶しておく。また、日本語の場合は特に、主語が省略されている文が多いので、機械的に翻訳する際に主語の補完が必要である。そこで非特許文献1に記載されている技術では、文毎に主語又はゼロ代名詞等の補完情報を付加する。これにより、当該情報が付加された文書を利用することによって正確に機械翻訳することが可能となる。文中で省略された単語、又は指示代名詞若しくはゼロ代名詞で表されている単語も、例えば文書を検索する場合の出現頻度の算出等の応用技術に利用することができる。

20

【非特許文献1】橋田浩一「大域文書修飾」人工知能学会全国大会（第11回）論文集 p . 62 - 63（1997）

30

【発明の開示】

【発明が解決しようとする課題】

【0012】

文章を書く時、又は発話する時の、その各文又は各発話夫々におけるユーザの注目対象（重点対象）は、会話や文章の文脈の流れに従って動的に変化する。つまり、会話や文章における単語への注目度合いを表す重みは、動的に変化する。よって、会話や文章に関連する情報を検索するサービスを実現するためには、文脈に応じた単語の重みの動的変化を追跡する必要がある。

【0013】

しかしながら、従来の文書検索サービスでは、検索のために入力された単語の出現頻度の高い文書を抽出し、抽出した文書から、当該単語を含む文又は段落を抽出して出力するため、当該単語のその文又は段落の文脈で動的に変わる重みについては考慮されずに検索される。したがって、出現頻度に基づく検索では、確かに検索のために入力された単語を含んではいるものの、文脈上当該単語がユーザが考えるように使用されていない場合があり、ユーザの検索目的を達成することができるとは限らない。各単語の文脈上の意味における各文での重み、即ち文脈上注目されているか否かについては特定できない。したがって、入力したキーワードをユーザの考える意味合い通りに使用した文又は段落を出力することはできない。

40

【0014】

また、非特許文献1の技術では、品詞情報等の文法に照らして識別が可能な情報を自動

50



的に解析し、指示代名詞又はゼロ代名詞等の補完、照応又は係り受けについての情報を文書に付加することができる。当該情報の付加により、参照されている名詞を出現頻度として利用することができるので、文又は段落等での単語間の関係は付加された情報により解析が可能である。しかしながら、各単語の文又は段落での注目されている度合い、即ち顕現性は、定量的に測ることはできない。

【 0 0 1 5 】

非特許文献 1 の技術は、自然文による質問に対して当該質問文で省略されている単語等を考慮してコンピュータに応答させる質問応答の実現へ応用が可能である。しかし、複数のユーザによる対話の文脈上の意味を定量的な値として算出し、第三者の発話としてユーザの対話の文脈に沿った発話を生成し、提示することを可能にするのは容易でない。

10

【 0 0 1 6 】

また、従来の文書検索サービスでは、文書中に出現する頻度が少ない場合でも文脈上深く関連する背景知識を表わすような単語を考慮して検索することはできなかつた。したがって、検索するユーザが意識しているが検索のために入力された単語としては現れていない単語を、同様に連想させる文又は段落を直接的に出力することはできなかつた。

【 0 0 1 7 】

本発明は斯かる事情に鑑みてなされたものであり、一又は複数の文からなる文単位毎に、その文単位での単語の顕現性を表わす重み値が夫々付与された重み付き単語群を対応付けて記憶しておき、検索のために受け付けた言葉についても、その言葉での重み値が付与された重み付き単語群を対応付け、重み付き単語群が類似する文単位を抽出して出力する構成とする。受け付けた言葉から、ユーザの意識にある先の言葉からの文脈が反映された意味を表わす情報を自動的に生成し、文脈の流れの中で意味が動的に変化する文書中の文単位の内から、受け付けた言葉から生成された情報が表わす文脈上の意味のまとまりが類似する文単位を直接的に検索することができる文単位検索方法、文単位検索装置、コンピュータを前記文単位検索装置として機能させるコンピュータプログラム、及び該コンピュータプログラムを記録したコンピュータ読み取り可能な記録媒体を提供することを目的とする。

20

【 0 0 1 8 】

本発明の目的は、文単位又は受け付ける言葉に対応付けられる重み付き単語群中の各単語の顕現性を表わす重み値を、後続の文単位又は言葉で出現する確率又は参照される確率として算出することにより、文脈の流れの中にある文単位又は言葉夫々で時系列に変化する単語の顕現性を定量的に表わして用いることができる文単位検索方法及び文書記憶装置を提供することにある。

30

【 0 0 1 9 】

また、本発明の目的は、関連する単語への関連度を定量的に算出し、各文単位又は言葉における各単語の顕現性に関連度を反映させることにより、ユーザから発せられる言葉又は筆記された文章には出現していない場合でも、ユーザが言葉を発しているとき又は筆記しているときに意識している単語を連想させる文単位をも効果的に検索することができる文単位検索方法及び文書記憶装置を提供することにある。

【課題を解決するための手段】

40

【 0 0 2 0 】

第 1 発明に係る文単位検索方法は、自然言語からなる複数の文書データが記憶されている文書集合を用い、該文書集合から取得した文書データを一又は複数の文からなる文単位に分別しておく一方、言葉を順次受け付け、受け付けた言葉に基づいて前記文書集合から分別してある文単位を検索する文単位検索方法において、文書データ中に連なる文単位夫々に、該文単位及び先行文脈に基づき求められる該文単位での顕現性を表わす重み値が付与された複数の単語からなる重み付き単語群を対応付けて予め記憶しておくステップと、言葉を受け付ける都度、該言葉に、該言葉及び先行文脈に基づき求められる該言葉での顕現性を表わす重み値が付与された複数の単語からなる重み付き単語群を対応付けるステップと、受け付けた言葉に対応付けた重み付き単語群と類似する重み付き単語群が対応付け

50

て記録されている文単位を、前記文書集合から抽出する類似文単位抽出ステップと、抽出した文単位を出力するステップとを含むことを特徴とする。

【0021】

第2発明に係る文単位検索方法は、前記類似文単位抽出ステップは、受け付けた言葉に対応付けた重み付き単語群の内の複数の単語の重み値の分布と、予め分別された文単位に対応付けられている重み付き単語群の内の複数の単語の重み値の分布とが、所定の条件を満たすか否かを判断するステップと、所定の条件を満たすと判断された重み付き単語群が対応付けられている文単位を抽出するステップとを含むことを特徴とする。

【0022】

第3発明に係る文単位検索方法は、前記類似文単位抽出ステップは、予め分別された文単位から、受け付けた言葉に対応付けた重み付き単語群と同一の単語を含む単語群が対応付けられた文単位を抽出するステップと、受け付けた言葉と抽出した文単位とで、対応付けられた単語群の内の同一の単語毎に重み値の差分を算出するステップと、抽出した文単位に、算出した差分が小さい順に優先順位を付与するステップとを含み、抽出した文単位を、優先順位に基づいて出力することを特徴とする。

10

【0023】

第4発明に係る文単位検索方法は、前記重み付き単語群を、各単語を1次元とし、単語毎に付与される重み値の大きさを各単語に対応する次元方向の要素として持つ多次元ベクトルとして算出するステップを含み、前記類似文単位抽出ステップは、分別した文単位毎に記憶してある前記多次元ベクトルと、受け付けた言葉に対応付けた前記多次元ベクトルとの距離を算出するステップと、文単位に、算出した距離が短い順に優先順位を付与するステップとを含み、付与された優先順位に従って出力することを特徴とする。

20

【0024】

第5発明に係る文単位検索方法は、文単位又は受け付けた言葉に重み付き単語群を対応付ける際、各単語が、前記文単位又は前記言葉よりも後続の文単位又は言葉に出現する又は参照される参照確率を算出する参照確率算出ステップを含み、算出した参照確率を各単語の重み値として付与することを特徴とする。

【0025】

第6発明に係る文単位検索方法は、前記参照確率算出ステップは、前記各単語が先行の文単位を含む複数の文単位に出現するパターン、又は前記単語を先行の文単位から参照するパターンを含む特徴パターンを特定するステップと、前記文書集合から取得された文書データ中で、前記特徴パターンと同一の特徴パターンが特定される単語が、後続の文単位で出現する又は参照される割合を算出するステップとを含み、算出した割合を参照確率とすることを特徴とする。

30

【0026】

第7発明に係る文単位検索方法は、前記文書集合から抽出される単語毎に、該単語が先行の文単位を含む複数の文単位に出現するパターン、又は前記単語を先行の文単位から参照するパターンを含む特徴パターンを特定する特定ステップと、特定した特徴パターンと同一の特徴パターンが特定される単語が、前記文書データ中で後続の文単位で出現したか又は参照されたかを判定する判定ステップと、特定した特徴パターンと、該特徴パターンで特定される単語に対して判定した結果との回帰分析を行って前記参照確率に対する前記特徴パターンの回帰係数を算出する回帰ステップとを含み、文単位に重み付き単語群を対応付けて記憶しておく際、又は受け付けた言葉に重み付き単語群を対応付ける際、前記参照確率算出ステップは、前記文単位又は言葉毎に、該文単位又は言葉での単語の特徴パターンを特定し、特定した特徴パターンに対する前記回帰係数を使用して参照確率を算出することを特徴とする。

40

【0027】

第8発明に係る文単位検索方法は、文単位に対しては、書き言葉からなる第1文書集合から取得された文書データ中で前記割合を算出し、受け付けた言葉に対しては、話し言葉からなる第2文書集合から取得された文書データ中で前記割合を算出することを特徴とす

50

る。

【0028】

第9発明に係る文単位検索方法は、書き言葉からなる第1文書集合及び話し言葉からなる第2文書集合夫々について、前記特定ステップ、前記判定ステップ及び前記回帰ステップを実行しておき、前記参照確率算出ステップは、前記文単位で特定した単語の特徴パターンに対しては、第1文書集合について実行した前記回帰ステップにより算出された回帰係数を使用して参照確率を算出し、前記受け付けた言葉で特定した単語の特徴パターンに対しては、第2文書集合について実行した前記回帰ステップで算出された回帰係数を使用して参照確率を算出することを特徴とする。

【0029】

第10発明に係る文単位検索方法は、前記特徴パターンは、前記単語を先行の文単位又は言葉から参照している場合の前記先行の文単位又は言葉から前記単語が含まれる文単位又は言葉までの、文単位又は言葉の数、前記単語が出現又は参照されている直近の先行の文単位又は言葉における前記単語の係り受け情報、前記単語が含まれる文単位又は言葉までに出現した又は参照された回数、前記単語が出現又は参照されている直近の先行の文単位又は言葉における前記単語の名詞区別、前記単語が出現又は参照されている直近の先行の文単位又は言葉中で前記単語が主題であるか否か、前記単語が出現又は参照されている直近の先行の文単位又は言葉中で前記単語が主語であるか否か、前記単語が含まれる文単位又は言葉における人称、及び、前記単語が含まれる文単位又は言葉における品詞情報、の内の一又は複数を含む情報で特定されることを特徴とする。

【0030】

第11発明に係る文単位検索方法は、前記特徴パターンは、前記単語を先行の文単位又は言葉から参照している場合の前記先行の文単位又は言葉から前記単語が含まれる文単位又は言葉までに対応する時間、前記単語が出現又は参照されている直近の先行の文単位又は言葉中で前記単語に対応する発話速度、及び、前記単語が出現又は参照されている直近の先行の文単位又は言葉中で前記単語に対応する音声の周波数の内の一又は複数を含む情報で特定されることを特徴とする。

【0031】

第12発明に係る文単位検索方法は、前記文章集合から抽出される単語の内の一の単語について、前記分別された文単位に対応付けられている重み付き単語群の内から、前記一の単語が含まれる単語群であり、且つ前記一の単語の重み値が所定値以上である単語群を抽出する第1ステップと、該第1ステップで抽出した単語群の各単語の重み値を単語毎に統合した値を、前記一の単語の各単語への関連度として付与した関連単語群を作成する第2ステップと、作成した関連単語群を前記一の単語に対応付けて記憶する第3ステップと、前記抽出された単語夫々について前記第1ステップ乃至第3ステップを予め実行するステップと、文単位毎又は受け付けた言葉毎に対応付けられた重み付き単語群の各単語の重み値夫々を、各単語に対応付けて記憶されている前記関連単語群の各単語の関連度を使用して付与し直す関連度付加ステップとを含むことを特徴とする。

【0032】

第13発明に係る文単位検索方法は、前記第2ステップは、前記抽出した単語群について、各単語群に含まれる各単語の重み値に、前記一の単語の重み値で重み付けした総和を算出するステップと、算出した総和を平均化するステップと、作成する関連単語群の各単語の前記関連度として、各単語の重み値の平均化された総和を付与するステップとを含むことを特徴とする。

【0033】

第14発明に係る文単位検索方法は、前記関連度付加ステップは、文単位毎又は受け付けた言葉毎に対応付けられた重み付き単語群の各単語について、各単語に対応付けて記憶されている前記関連単語群に含まれる各単語の関連度を、前記重み付き単語群の各単語の重み値に乗算するステップと、乗算結果に基づいて前記重み付き単語群の各単語の重み値として付与し直すステップとを含むことを特徴とする。

## 【 0 0 3 4 】

第 1 5 発明に係る文単位検索方法は、各単語夫々についての前記関連単語群を、各単語を 1 次元とし、単語毎に付与される関連度の大きさを各単語に対応する次元方向の要素として持つ多次元の関連度ベクトルとして算出するステップとを含み、前記関連度付加ステップは、分別した文単位毎に記憶してある前記多次元ベクトルを、各単語の関連度ベクトルの列によって変換することを特徴とする。

## 【 0 0 3 5 】

第 1 6 発明に係る文単位検索方法は、自然言語からなる複数の文書データが記憶されている文書集合を用い、言葉を受け付け、受け付けた言葉に基づいて前記文書集合を検索する文単位検索方法において、前記文書集合から得られる文書データを一又は複数の文からなる文単位に分別しておくステップ、分別した文単位毎に、該文単位に出現する単語、又は、文書データ中の先行の文単位から参照する単語を抽出するステップ、前記文単位に対して抽出した単語毎に、各文単位における特徴を特定して記憶しておくステップ、分別した文単位毎に、該文単位に対して抽出した単語が該文単位及び先行の文単位で出現する場合の前記特徴の組み合わせのパターン、又は先行の文単位から参照する場合の参照のパターンを含む特徴パターンを特定するステップ、特定した特徴パターンと、該特徴パターンで特定された単語が後続の文単位で出現又は参照されたか否かとを記憶しておくステップ、前記文書集合から得られる文書中の文単位全体に対し、一の特徴パターンで特定される単語が後続の文単位で出現又は参照される参照確率の回帰分析を行って特徴パターンに対応する回帰係数を得る回帰学習を実行するステップ、分別した文単位毎に、文書データ中 10  
20  
30  
40  
50  
60  
70  
80  
90  
100  
110  
120  
130  
140  
150  
160  
170  
180  
190  
200  
210  
220  
230  
240  
250  
260  
270  
280  
290  
300  
310  
320  
330  
340  
350  
360  
370  
380  
390  
400  
410  
420  
430  
440  
450  
460  
470  
480  
490  
500  
510  
520  
530  
540  
550  
560  
570  
580  
590  
600  
610  
620  
630  
640  
650  
660  
670  
680  
690  
700  
710  
720  
730  
740  
750  
760  
770  
780  
790  
800  
810  
820  
830  
840  
850  
860  
870  
880  
890  
900  
910  
920  
930  
940  
950  
960  
970  
980  
990  
1000  
1010  
1020  
1030  
1040  
1050  
1060  
1070  
1080  
1090  
1100  
1110  
1120  
1130  
1140  
1150  
1160  
1170  
1180  
1190  
1200  
1210  
1220  
1230  
1240  
1250  
1260  
1270  
1280  
1290  
1300  
1310  
1320  
1330  
1340  
1350  
1360  
1370  
1380  
1390  
1400  
1410  
1420  
1430  
1440  
1450  
1460  
1470  
1480  
1490  
1500  
1510  
1520  
1530  
1540  
1550  
1560  
1570  
1580  
1590  
1600  
1610  
1620  
1630  
1640  
1650  
1660  
1670  
1680  
1690  
1700  
1710  
1720  
1730  
1740  
1750  
1760  
1770  
1780  
1790  
1800  
1810  
1820  
1830  
1840  
1850  
1860  
1870  
1880  
1890  
1900  
1910  
1920  
1930  
1940  
1950  
1960  
1970  
1980  
1990  
2000  
2010  
2020  
2030  
2040  
2050  
2060  
2070  
2080  
2090  
2100  
2110  
2120  
2130  
2140  
2150  
2160  
2170  
2180  
2190  
2200  
2210  
2220  
2230  
2240  
2250  
2260  
2270  
2280  
2290  
2300  
2310  
2320  
2330  
2340  
2350  
2360  
2370  
2380  
2390  
2400  
2410  
2420  
2430  
2440  
2450  
2460  
2470  
2480  
2490  
2500  
2510  
2520  
2530  
2540  
2550  
2560  
2570  
2580  
2590  
2600  
2610  
2620  
2630  
2640  
2650  
2660  
2670  
2680  
2690  
2700  
2710  
2720  
2730  
2740  
2750  
2760  
2770  
2780  
2790  
2800  
2810  
2820  
2830  
2840  
2850  
2860  
2870  
2880  
2890  
2900  
2910  
2920  
2930  
2940  
2950  
2960  
2970  
2980  
2990  
3000  
3010  
3020  
3030  
3040  
3050  
3060  
3070  
3080  
3090  
3100  
3110  
3120  
3130  
3140  
3150  
3160  
3170  
3180  
3190  
3200  
3210  
3220  
3230  
3240  
3250  
3260  
3270  
3280  
3290  
3300  
3310  
3320  
3330  
3340  
3350  
3360  
3370  
3380  
3390  
3400  
3410  
3420  
3430  
3440  
3450  
3460  
3470  
3480  
3490  
3500  
3510  
3520  
3530  
3540  
3550  
3560  
3570  
3580  
3590  
3600  
3610  
3620  
3630  
3640  
3650  
3660  
3670  
3680  
3690  
3700  
3710  
3720  
3730  
3740  
3750  
3760  
3770  
3780  
3790  
3800  
3810  
3820  
3830  
3840  
3850  
3860  
3870  
3880  
3890  
3900  
3910  
3920  
3930  
3940  
3950  
3960  
3970  
3980  
3990  
4000  
4010  
4020  
4030  
4040  
4050  
4060  
4070  
4080  
4090  
4100  
4110  
4120  
4130  
4140  
4150  
4160  
4170  
4180  
4190  
4200  
4210  
4220  
4230  
4240  
4250  
4260  
4270  
4280  
4290  
4300  
4310  
4320  
4330  
4340  
4350  
4360  
4370  
4380  
4390  
4400  
4410  
4420  
4430  
4440  
4450  
4460  
4470  
4480  
4490  
4500  
4510  
4520  
4530  
4540  
4550  
4560  
4570  
4580  
4590  
4600  
4610  
4620  
4630  
4640  
4650  
4660  
4670  
4680  
4690  
4700  
4710  
4720  
4730  
4740  
4750  
4760  
4770  
4780  
4790  
4800  
4810  
4820  
4830  
4840  
4850  
4860  
4870  
4880  
4890  
4900  
4910  
4920  
4930  
4940  
4950  
4960  
4970  
4980  
4990  
5000  
5010  
5020  
5030  
5040  
5050  
5060  
5070  
5080  
5090  
5100  
5110  
5120  
5130  
5140  
5150  
5160  
5170  
5180  
5190  
5200  
5210  
5220  
5230  
5240  
5250  
5260  
5270  
5280  
5290  
5300  
5310  
5320  
5330  
5340  
5350  
5360  
5370  
5380  
5390  
5400  
5410  
5420  
5430  
5440  
5450  
5460  
5470  
5480  
5490  
5500  
5510  
5520  
5530  
5540  
5550  
5560  
5570  
5580  
5590  
5600  
5610  
5620  
5630  
5640  
5650  
5660  
5670  
5680  
5690  
5700  
5710  
5720  
5730  
5740  
5750  
5760  
5770  
5780  
5790  
5800  
5810  
5820  
5830  
5840  
5850  
5860  
5870  
5880  
5890  
5900  
5910  
5920  
5930  
5940  
5950  
5960  
5970  
5980  
5990  
6000  
6010  
6020  
6030  
6040  
6050  
6060  
6070  
6080  
6090  
6100  
6110  
6120  
6130  
6140  
6150  
6160  
6170  
6180  
6190  
6200  
6210  
6220  
6230  
6240  
6250  
6260  
6270  
6280  
6290  
6300  
6310  
6320  
6330  
6340  
6350  
6360  
6370  
6380  
6390  
6400  
6410  
6420  
6430  
6440  
6450  
6460  
6470  
6480  
6490  
6500  
6510  
6520  
6530  
6540  
6550  
6560  
6570  
6580  
6590  
6600  
6610  
6620  
6630  
6640  
6650  
6660  
6670  
6680  
6690  
6700  
6710  
6720  
6730  
6740  
6750  
6760  
6770  
6780  
6790  
6800  
6810  
6820  
6830  
6840  
6850  
6860  
6870  
6880  
6890  
6900  
6910  
6920  
6930  
6940  
6950  
6960  
6970  
6980  
6990  
7000  
7010  
7020  
7030  
7040  
7050  
7060  
7070  
7080  
7090  
7100  
7110  
7120  
7130  
7140  
7150  
7160  
7170  
7180  
7190  
7200  
7210  
7220  
7230  
7240  
7250  
7260  
7270  
7280  
7290  
7300  
7310  
7320  
7330  
7340  
7350  
7360  
7370  
7380  
7390  
7400  
7410  
7420  
7430  
7440  
7450  
7460  
7470  
7480  
7490  
7500  
7510  
7520  
7530  
7540  
7550  
7560  
7570  
7580  
7590  
7600  
7610  
7620  
7630  
7640  
7650  
7660  
7670  
7680  
7690  
7700  
7710  
7720  
7730  
7740  
7750  
7760  
7770  
7780  
7790  
7800  
7810  
7820  
7830  
7840  
7850  
7860  
7870  
7880  
7890  
7900  
7910  
7920  
7930  
7940  
7950  
7960  
7970  
7980  
7990  
8000  
8010  
8020  
8030  
8040  
8050  
8060  
8070  
8080  
8090  
8100  
8110  
8120  
8130  
8140  
8150  
8160  
8170  
8180  
8190  
8200  
8210  
8220  
8230  
8240  
8250  
8260  
8270  
8280  
8290  
8300  
8310  
8320  
8330  
8340  
8350  
8360  
8370  
8380  
8390  
8400  
8410  
8420  
8430  
8440  
8450  
8460  
8470  
8480  
8490  
8500  
8510  
8520  
8530  
8540  
8550  
8560  
8570  
8580  
8590  
8600  
8610  
8620  
8630  
8640  
8650  
8660  
8670  
8680  
8690  
8700  
8710  
8720  
8730  
8740  
8750  
8760  
8770  
8780  
8790  
8800  
8810  
8820  
8830  
8840  
8850  
8860  
8870  
8880  
8890  
8900  
8910  
8920  
8930  
8940  
8950  
8960  
8970  
8980  
8990  
9000  
9010  
9020  
9030  
9040  
9050  
9060  
9070  
9080  
9090  
9100  
9110  
9120  
9130  
9140  
9150  
9160  
9170  
9180  
9190  
9200  
9210  
9220  
9230  
9240  
9250  
9260  
9270  
9280  
9290  
9300  
9310  
9320  
9330  
9340  
9350  
9360  
9370  
9380  
9390  
9400  
9410  
9420  
9430  
9440  
9450  
9460  
9470  
9480  
9490  
9500  
9510  
9520  
9530  
9540  
9550  
9560  
9570  
9580  
9590  
9600  
9610  
9620  
9630  
9640  
9650  
9660  
9670  
9680  
9690  
9700  
9710  
9720  
9730  
9740  
9750  
9760  
9770  
9780  
9790  
9800  
9810  
9820  
9830  
9840  
9850  
9860  
9870  
9880  
9890  
9900  
9910  
9920  
9930  
9940  
9950  
9960  
9970  
9980  
9990  
10000

## 【 0 0 3 6 】

第 1 7 発明に係る文単位検索装置は、自然言語からなる複数の文書データが記憶されている文書集合から文書データを取得する手段と、言葉を順次受け付ける手段とを備え、受け付けた言葉に基づいて前記文書集合を検索する文単位検索装置において、取得した文書データを一又は複数の文からなる文単位に分別する手段と、取得した文書データ中に連なる文単位夫々に、該文単位及び先行文脈に基づき求められる該文単位での顕現性を表わす重み値が付与された複数の単語からなる重み付き単語群を対応付けて記憶する手段と、言葉を受け付けた場合に受け付けた順に記憶する手段と、新たに言葉を受け付ける都度、該言葉に、該言葉及び該先行文脈に基づき求められる該言葉での顕現性を表わす重み値が付与された複数の単語からなる重み付き単語群を対応付ける手段と、予め分別された文単位から、受け付けた言葉に対応付けた重み付き単語群と類似する重み付き単語群が対応付けて記録されている文単位を抽出する手段と、抽出した文単位を出力する手段とを備えることを特徴とする。

## 【 0 0 3 7 】

第 1 8 発明に係るコンピュータプログラムは、自然言語からなる複数の文書データが記

10

20

30

40

50

憶されている文書集合から、文書データを取得することが可能であるコンピュータを、言葉を順次受け付ける手段と、受け付けた言葉に基づいて前記文書集合を検索する手段として機能させることができるコンピュータプログラムにおいて、取得した文書データを一又は複数の文からなる文単位に分別する手段、取得した文書データ中に連なる文単位夫々に、該文単位及び先行文脈に基づき求められる該文単位での顕現性を表わす重み値が付与された複数の単語からなる重み付き単語群を対応付けて記憶する手段、言葉を受け付けた場合に受け付けた順に記憶する手段、新たに言葉を受け付ける都度、該言葉に、該言葉及び先行文脈に基づき求められる該言葉での顕現性を表わす重み値が付与された複数の単語からなる重み付き単語群を対応付ける手段、及び、予め分別された文単位から、受け付けた言葉に対応付けた重み付き単語群と類似する重み付き単語群が対応付けて記録されている文単位を抽出する手段として機能させることを特徴とする。

10

【0038】

第19発明に係るコンピュータで読み取り可能な記録媒体には、第18発明のコンピュータプログラムが記録されていることを特徴とする。

【0039】

第20発明に係る文書記憶装置は、自然言語からなる複数の文書データを記憶する手段と、記憶した文書データを、文書データの先頭から順に一又は複数の文からなる文単位に分別する手段とを備え、分別した文単位毎に、該文単位に出現する単語又は先行する文単位から参照する単語が抽出してあり、分別した文単位毎に抽出した単語が記憶してある文書記憶装置において、文書データ中に連なる文単位毎に、複数の単語が、該文単位よりも後続の文単位に出現するか又は参照される参照確率を算出する手段と、前記文単位夫々に、該文単位での顕現性を表わす重み値として前記参照確率が付与された前記複数の単語からなる重み付き単語群を対応付けて記憶する手段とを備えることを特徴とする。

20

【0040】

第21発明に係る文書記憶装置は、抽出されてある単語の内の一の単語について、文単位夫々に対応付けられている重み付き単語群の内から、前記一の単語が含まれる単語群であり、且つ前記一の単語の重み値が所定値以上である単語群を抽出する抽出手段と、該抽出手段が抽出した単語群の各単語の重み値を単語毎に統合した値を、前記一の単語の各単語への関連度として付与した関連単語群を作成する作成手段と、作成した関連単語群を前記一の単語に対応付けて記憶する記憶手段とを備え、前記抽出されてある単語夫々について前記抽出手段、前記作成手段及び前記記憶手段の処理を実行するようにしてあり、各単語に対応付けて夫々の関連単語群を記憶するようにしてあることを特徴とする。

30

【0041】

第1発明、第17発明、第18発明及び第19発明では、自然言語からなる文書データが記録された文書集合から文書データが取得され、取得された文書データは更に一又は複数の文である文単位に分別される。文単位毎に、文書集合中で出現する各単語についてその文単位での重み値が付与され、重み値が付与された単語の重み付き単語群が文単位に対応付けて記憶される。言葉を受け付けた場合、受け付けた言葉についてもその言葉での重み値が付与された単語の重み付き単語群が対応付けられる。予め分別されている文単位から、受け付けた言葉に対応付けられた重み付き単語群と類似する重み付き単語群が対応付けられている文単位が抽出され、出力される。

40

【0042】

第2発明では、第1発明において類似する重み付き単語群が対応付けられている文単位を抽出する際、予め文単位に対応付けて記憶されている重み付き単語群の内の複数の単語の重み値の分布が、受け付けた言葉に対応付けられた重み付き単語群の内の複数の単語の重み値の分布と所定の条件を満たすか否かの判断により類似するか否かが判定され、類似すると判定された重み付き単語群が対応付けられている文単位が抽出される。

【0043】

第3発明では、第1発明又は第2発明において類似する重み付き単語群が対応付けられている文単位を抽出する際、重み付き単語群に同一の単語が含まれる文単位が抽出され、

50

その同一の単語に付与されている重み値の差分が小さい順に優先順位が付与される。

【 0 0 4 4 】

第4発明では、第1発明における重み付き単語群は、各単語を1次元とし、単語毎に付与される重み値の大きさを各単語に対応する次元方向の要素として持つ多次元ベクトルとして得られる。重み付き単語群が類似するか否かの判定を、重み付き単語群同士、即ち多次元ベクトル間の距離が短いかなかで判定される。抽出された文単位は、多次元ベクトル間の距離が短い順、即ち重み付き単語群同士が類似する順に出力される。

【 0 0 4 5 】

第5発明では、第1発明乃至第4発明において各単語に付与される重み値として、各単語が夫々、後続の文単位又は言葉に出現する又は参照される参照確率が算出されて付与される。

10

【 0 0 4 6 】

第6発明では、第5発明において算出される参照確率は、各単語に対して特定される先行の文単位から各文単位に至るまでの出現のパターン、又は先行の文単位からの参照のパターンを含む特徴パターンと同一の特徴パターンが特定される単語が、文書集合中で後続の文単位でさらに出現する又は参照される割合として算出される。

【 0 0 4 7 】

第7発明では、文書集合から抽出される各単語に対し特定される特徴パターンと、その特徴パターンが特定される単語が文書集合中の文書中の後続の文単位で出現したか又は参照されたかの判定結果とが回帰分析され、単語が後続の文単位で出現又は参照される参照確率に対する特徴パターンの回帰係数が算出される。第5発明において算出される参照確率は、単語毎に夫々の特徴パターンが特定され、その特徴パターンと回帰係数とから算出される。

20

【 0 0 4 8 】

第8発明及び第9発明では、文書集合が書き言葉からなる第1文書集合と、話し言葉からなる第2文書集合とに分けられて用いられる。文単位に対応付けられる重み付き単語群の各単語へ付与する参照確率は、第1文書集合に基づいて算出され、受け付けた言葉に対応付けられる重み付き単語群の各言葉へ付与する参照確率は、第2文書集合に基づいて算出される。

【 0 0 4 9 】

30

第10発明では、第6発明乃至第9発明において参照確率を算出する際に、各単語の特徴パターンを特定するための特徴として、先行の文単位又は言葉で出現又は参照している場合の現在の文単位又は言葉に至るまでの数、出現又は参照した場合の単語の係り受け情報、出現した回数又は参照された回数、単語の名詞区別、単語が主題であるか、単語が主語であるか、単語の人称、単語の品詞情報等の情報が定量的に扱われる。

【 0 0 5 0 】

第11発明では、第6発明乃至第10発明において参照確率を算出する際に、各単語の特徴パターンを特定するための特徴として、先行の文単位又は言葉で出現又は参照している場合に先行の文単位又は言葉からの時間、出現又は参照した場合のその単語に相当する音声の発話速度、音声の周波数の高低の情報が定量的に扱われる。

40

【 0 0 5 1 】

第12発明では、第1発明乃至第11発明において、文書集合から抽出される単語の内の一の単語について、その単語の重み値が所定値以上の重み付き単語群が抽出される。その一の単語について抽出された複数の重み付き単語群の各単語の重み値を単語毎に統合した一の重み付き単語群が関連単語群として作成される。作成された関連単語群の各単語の関連度は、一の単語に所定値以上の重み値が付与されている場合の各単語の重み値への関連の深さを表わしている。文書集合から抽出される単語夫々に対して関連単語群が生成され記憶される。各文単位又は言葉に対応付けられた重み付き単語群の各単語の重み値が、夫々の単語に対応付けられた関連単語群の各単語の関連度を使用して付与し直される。

【 0 0 5 2 】

50

第13発明では、第12発明において一の単語に対する関連単語群が作成される際、一の単語の重み値が所定値以上である重み付き単語群として抽出された単語群が、その重み付き単語群での前記一の単語に対する重み値によって重み付けされた総和が算出される。総和は平均化され、各単語について平均化された重み値の総和が関連単語群の各単語の関連度として付与される。

【0053】

第14発明では、前記12発明又は第13発明で記憶される関連単語群の各単語の関連度が、文単位毎又は受け付けた言葉毎に対応付けられた重み付き単語群の各単語の重み値に乗算され、乗算結果が重み付き単語群の各単語の重み値として付与し直される。重み付き単語群の内の一の単語に注目した場合、一の単語に対応付けられた関連単語群の各単語の関連度が使用される。重み付き単語群の内の一の単語以外の各単語の重み値と、前記一の単語に対応付けられた関連単語群の各単語の関連度とが乗算されることにより、関連度の高い他の単語の重み値からの前記一の単語の重み値への影響が加味される。

10

【0054】

第15発明では、第12発明乃至第14発明における関連単語群は、各単語を1次元とし、単語毎に付与される関連度の大きさを各単語に対応する次元方向の要素として持つ多次元の関連度ベクトルとして得られる。各文単位又は言葉に対応付けられた多次元ベクトルは、各単語に対する関連語ベクトルの列からなる行列で変換される。即ち、多次元ベクトルは単語の各1次元間の距離が関連度が高い単語の次元間ほど距離が短い斜交座標系における多次元ベクトルで表現される。したがって、重み付き単語群を表現する多次元ベクトルは、それに含まれる単語と関連度が高い単語軸方向に回転され、関連度が高い単語を含む多次元ベクトル間の距離はより短くなる。

20

【0055】

第16発明では、文書集合から取得された文書データを更に分別した文単位毎に、文単位又は先行の文単位から参照する単語が抽出され、各単語に対して各文単位における特徴が特定され、先行の文単位から各文単位に至るまでの特徴の組み合わせのパターン、又は各単語の先行の文単位からの参照のパターンを含む特徴パターンが特定される。特定された特徴パターンによる参照確率の回帰学習に基づいて、抽出された各単語の参照確率が算出され、重み付き単語群として予め文単位毎に記憶される。受け付けた言葉に対しても先行の言葉に基づいた特徴パターンが特定されて各単語の参照確率が算出され、重み付き単語群が対応付けられる。予め記憶してある文単位は、受け付けた言葉の重み付き単語群と同一の単語の参照確率の差分が小さい順に優先順位が付与されて出力される。

30

【0056】

第20発明では、文書集合から取得された文書データを更に分別した文単位毎に、その文単位での単語の重みが付与された重み付き単語群が対応付けられて記憶される。

【0057】

第21発明では、第12発明で文書から抽出されてある単語夫々について作成された関連単語群が記憶される。

【発明の効果】

【0058】

本発明による場合、文書集合から取得した文書データ中の一又は複数の文からなる文単位毎に、複数の単語夫々の当該文単位での重み値を付与した重み付き単語群が対応付けられて記憶される。重み値付き単語群は、各文単位での各単語の重み値の組であり、文単位毎の意味のまとまりを示す情報として推定することができる。各重み値に先行の文単位から続く文脈が反映された値が付与されていることにより、分別された連なる文単位中の各文単位での重み付き単語群は、文書全体での意味のまとまりと異なり、文書中にある先行の文から続く文脈の流れの中で、動的に時系列的に変化していく意味のまとまりとして捉えることができる。検索のために入力される言葉での重み値が付与された重み付き単語群と類似する重み付き単語群が対応付けられる文単位が抽出されることにより、文書全体ではなく、単語の顕現性、即ち意味のまとまりが類似する文単位を直接的に検索することが

40

50

できる。

【 0 0 5 9 】

また、重み付き単語群が類似するか否かは、受け付けた言葉の重み付き単語群の内の複数の単語の重み値の分布と、予め記憶してある重み付き単語群の内の複数の単語の重み値の分布とを比較した場合に、分布同士が類似であると判断できる所定の条件を満たすとき、記憶してある重み付き単語群が受け付けた言葉の重み付き単語群と類似するといえることができる。例えば、重み付き単語群同士が類似しているといえる条件とした場合、重み付き単語群が類似しているといえることができる。つまり、一方の重み付き単語群において一の単語の重み値の他の単語の重み値に対する比率が、他方の重み付き単語群における一の単語の重み値の他の単語の重み値に対する比率にも保存される場合、それらの重み付き単語群同士は類似していると判断することができる。また、所定の条件を、例えば、一又は複数の単語に注目した場合にその単語の重み値がいずれも所定値以上であるか否かに設定することで判断することもできる。また、受け付けた言葉に対応付けた重み付き単語群と、予め分別された文単位に対応付けられている重み付き単語群と比較した場合に、同一の単語の重み値の差が小さいか否かにより類似するか否かを判断することもできる。

10

【 0 0 6 0 】

また、重み付き単語群を、各単語を1次元として、各単語の文単位又は言葉での重み値を各次元成分に対する要素として持つ多次元ベクトルとして表現することにより、文単位又は言葉毎の意味のまとまりを定量的なベクトルとして扱うことができる。また、文単位又は言葉毎の意味のまとまりを定量的な多次元ベクトルとして扱うことにより、ベクトル演算が可能なコンピュータを利用して、受け付けた言葉に対応付けられたベクトルと記憶してある文単位毎に対応付けられたベクトルとの距離を算出することによって類似する文単位を直接的に抽出することができる。さらに、多次元ベクトルとして表現することによって、受け付けた言葉、又は予め分別された文単位の多次元ベクトルが満たす条件を、多次元空間上のどの空間に相当するか否かによって設定することができ、類似する文単位を直接的に抽出することができる。

20

【 0 0 6 1 】

なお、ここでいう文書集合は、いわゆる書き言葉からなる文書データの集合に限らない。したがって、それらを分別した文単位も書き言葉からなる文単位とは限らない。文書データは既に記憶されてあるデータを意味してリアルタイムに受け付ける言葉と区別するものであり、話し言葉による対話が順に書下された文書データでもよい。

30

【 0 0 6 2 】

また、受け付ける言葉は、検索の目的で入力される単語、文章等に限らず、例えばユーザ同士の対話中の音声を含む各発話でもよい。各発話での重み値が付与された重み付き単語群に基づいて文単位を抽出するので、対話中で発話毎に意味が動的に、時系列的に変化していくことを考慮した意味のまとまりを発話毎に推定することができる。したがって、各発話に対して推定される意味のまとまりに類似する文単位を抽出して提示することが可能になる。

【 0 0 6 3 】

さらに、本発明による場合、重み付き単語群の各単語の重み値を、後続の文単位又は言葉でも出現又は参照される参照確率として付与することにより、各単語の重み値を注目されている度合い、即ち顕現性を示す定量的な値で表わすことができる。文脈上のその文単位において重要な注目されている単語は、継続して出現又は参照される確率が高いと考えられる。したがって、参照確率はその文単位における各単語の注目されている度合い、即ち顕現性を示すといえることができる。

40

【 0 0 6 4 】

また、各文単位で実際に出現することなしに指示代名詞又はゼロ代名詞で表わされる単語、又は指示代名詞又はゼロ代名詞でも表わされていない単語であっても、文単位又は言葉に実際に出現していない単語であっても後続の文単位又は言葉で出現又は参照される単

50



語は、その文単位又は言葉での顕現性が高いと考えられる。各文単位を基準とした先行の複数の文単位での単語の特徴パターンに基づいて参照確率を算出するので、実際に出現していない単語であっても、顕現性の高さをより正しく定量的に表わすことができる。

【0065】

さらに、言葉を音声で受け付けた場合は、言葉が発声されたときの声の特徴、即ち話す速度、声調からも、その言葉に含まれる単語がその言葉で重みを持っているのか否かを定量的に特徴づけて各単語の顕現性の高さを表わすことができる。

【0066】

さらに、本発明による場合、検索結果として出力する文単位が書き言葉である場合は、書き言葉からなる文書集合に基づいて参照確率を算出し、受け付けた言葉が話し言葉である場合は、話し言葉からなる文書集合に基づいて参照確率を学習、算出する。これにより、書き言葉と話し言葉とで異なる特徴を踏まえて、より意味合いが似た文単位を出力することができる。

10

【0067】

また、本発明による場合、単語毎に各単語からの関連度を定量的に算出して記憶しておく。重み付き単語群の内の各単語の重み値を、他の単語の重み値と、各単語からの前記一への単語の関連度とに基づいて算出し直す。これにより、一の単語の重み値に対し、他の単語の内の一の単語に対する関連度が高い単語の重み値の影響を反映させることができる。つまり、一の単語に対する関連度が高い単語の重み値が高い場合は、一の単語の重み値が高くなることを再現することができる。

20

【0068】

一の単語に対する関連語群を関連度ベクトルとして表現し、重み付き単語群を多次元ベクトルで表現した場合に各単語に対する関連度ベクトルの列からなる行列で多次元ベクトルを変換することにより、関連度の強い単語を含む重み付き単語群を表現する多次元ベクトル間の距離が短くなる。

【0069】

これにより、重み付き単語群の内の一の単語以外の単語の内、前記一の単語への関連度が高い単語の重み値の影響を、前記一の単語の重み値に反映することができる。各文単位又は言葉での各単語の顕現性に関連度を反映させて、受け付けた言葉に表れていない場合であってもユーザに意識されている単語を連想させる文単位を効果的に検索することができる等の優れた効果を奏する。

30

【図面の簡単な説明】

【0070】

【図1】本発明に係る文単位検索方法の概要を示す説明図である。

【図2】実施の形態1における文単位検索装置を用いた検索システムの構成を示すブロック図である。

【図3】実施の形態1における文単位検索装置のCPUが、取得した文書データに対する形態素解析及び統語解析処理の解析結果からタグ付け及び単語抽出を行い記憶する処理手順を示すフローチャートである。

【図4】実施の形態1における文書記憶手段で記憶される文書データの内容の一例を示す説明図である。

40

【図5】実施の形態1における文単位検索装置のCPUが、形態素解析及び統語解析した結果を付与して文書記憶手段に記憶させる文書データの一例を示す説明図である。

【図6】実施の形態1における文単位検索装置のCPUが取得した全文書データから抽出した単語のリストの例を示す説明図である。

【図7】実施の形態1における文単位検索装置のCPUが、文書記憶手段で記憶しているタグ付け済み文書データからサンプルを抽出し、回帰分析を行って参照確率を算出するための回帰式を推定する処理手順を示すフローチャートである。

【図8】実施の形態1における文書記憶手段で記憶された文書データ中の文で特定される特徴パターンの例を示す説明図である。

50

【図 9】実施の形態 1 における文単位検索装置の CPU が、文書記憶手段で記憶しているタグ付け済みの文書データの文毎に単語の参照確率を算出し、記憶する処理手順を示すフローチャートである。

【図 10】実施の形態 1 における文単位検索装置の CPU が、文書記憶手段で記憶しているタグ付け済みの文書データの文毎に単語の参照確率を算出し、記憶する処理手順を示すフローチャートである。

【図 11】実施の形態 1 における文単位検索装置の CPU が、文書データに示される文書を文毎に分別した一例を示す説明図である。

【図 12】実施の形態 1 における文単位検索装置の CPU が、参照確率を算出した結果を付与して文書記憶手段に記憶させる文書データの一例を示す説明図である。

10

【図 13】実施の形態 1 における文単位検索装置の CPU が、文単位毎に算出した重み付き単語群を索引付けして記憶した場合のデータベースの内容例を示す説明図である。

【図 14】文単位検索装置の CPU により文毎に記憶される単語及び該単語に対して算出された参照確率の組が、文が続くにつれてどのように変化するかを示す説明図である。

【図 15】実施の形態 1 における文単位検索装置及び受付装置の検索処理の処理手順を示すフローチャートである。

【図 16】実施の形態 1 における文単位検索装置及び受付装置の検索処理の処理手順を示すフローチャートである。

【図 17】実施の形態 1 における文単位検索装置及び受付装置の検索処理の処理手順を示すフローチャートである。

20

【図 18】実施の形態 1 における文単位検索装置の CPU が、受付装置から受信したテキストデータに対して特定した特徴パターンの例を示す説明図である。

【図 19】実施の形態 2 における文単位検索装置及び受付装置の検索処理の処理手順を示すフローチャートである。

【図 20】実施の形態 3 における本発明の検索方法に関わる、一の単語と関連の深い単語の顕現性の影響の概要を示す説明図である。

【図 21】実施の形態 3 における文単位検索装置の CPU が関連語群を作成する処理手順を示すフローチャートである。

【図 22】実施の形態 3 における文単位検索装置の CPU が関連語群を作成する処理手順を示すフローチャートである。

30

【図 23】実施の形態 3 における文単位検索装置の CPU によって関連語群が作成される場合の、各処理の過程での重み付き単語群の例を示す説明図である。

【図 24】実施の形態 3 における文単位検索装置の CPU が、各文単位に対応付けられて記憶されている重み付き単語群の各単語の重み値を算出し直す処理手順を示すフローチャートである。

【図 25】実施の形態 3 における文単位検索装置の CPU が、各文単位に対応付けられて記憶されている重み付き単語群の各単語の重み値を算出し直す処理手順の詳細を示すフローチャートである。

【図 26】実施の形態 3 における文単位検索装置の CPU によって算出された各単語の顕現性を表わす重み値の内容例を示す説明図である。

40

【図 27】実施の形態 3 における文単位検索装置及び受付装置の検索処理の処理手順を示すフローチャートである。

【図 28】実施の形態 3 における文単位検索装置及び受付装置の検索処理の処理手順を示すフローチャートである。

【図 29】本発明の文単位検索方法を文単位検索装置で実施する場合の構成を示すブロック図である。

【符号の説明】

【0071】

1 文単位検索装置

11 CPU

50

- 1 3 記憶手段
- 1 5 通信手段
- 1 6 文書集合接続手段
- 1 7 補助記憶手段
- 1 8 可搬型記録媒体
- 1 P 制御プログラム
- 2 文書記憶手段
- 4 受付装置

【発明を実施するための最良の形態】

【0072】

以下本発明をその実施の形態を示す図面に基づき具体的に説明する。

【0073】

図1は、本発明に係る文単位検索方法の概要を示す説明図である。図1中の100は、複数の文書データが記憶されている文書集合を表わしており、文書集合100から取得される一の文書101は、一又は複数の文からなる文単位 $S_1, \dots, S_i, S_{i+1}, \dots$ で構成されている。文単位 $S_1, \dots, S_i, S_{i+1}, \dots$ は、文書101の先頭から順に文脈の流れに沿い、時系列的に変遷する意味合いを有して連なっている。図1中の200は、ユーザAとユーザBとの会話を表わしている。ユーザAとユーザBの会話200は、上から下へ時系列に連なるユーザA及びユーザBからの発話 $U_{j-3}, \dots, U_j$ の集合である。会話は、発話 $U_{j-3}, U_{j-2}, U_{j-1}, U_j$ の順になされている。なお、ユーザAとユーザB

10

20

【0074】

本発明に係る文単位検索方法は、文単位又は言葉をユーザが筆記又は発話した時点での各単語への注目度合いを定量的な重み値として表わして各単語に付与し、時系列に連続する文単位又は言葉毎に変遷していく各単語への注目度合いを反映した重み付き単語群を各文単位における文脈上の意味合いを表わす指標として用いることにより、同様の文脈上の意味合いを有する文単位を直接的に検索し、出力することを目的としている。

【0075】

図1の説明図で示す例での会話200は、ユーザAとユーザBとの間でなされている京都への旅行についての会話である。会話200中の発話 $U_{j-3}$ では「京都」「旅行」が現れ、文脈の流れは「京都の旅行」である。発話 $U_{j-2}$ では、「京都」「旅行」は現れていないが「“京都への旅行の”時期」についての発話であり、「京都」「旅行」「時期」について注目がされている。 $U_{j-1}$ では「暑い」が現れている。 $U_{j-1}$ では「京都」「旅行」は現れていないが、「“京都は”暑い」のであり、依然「京都」は文脈上の意味に対して重みを持っている。さらにユーザAとユーザBの間では、 $U_{j-1}$ の発話の時点では、「旅行」よりも「京都」及び「時期」が注目されており、ユーザAとユーザBとは文脈上の意味合いが変遷していることを共通して認識できるはずである。さらに、発話 $U_j$ の中で「有名」「祭」が現れている。この $U_j$ の発話の時点だけを考えれば、「京都」「旅行」「時期」「暑い」という単語は現れていない。しかし、少なくともユーザAにとっては、発話 $U_j$ は文脈上「夏」の「京都」の「祭」についての意味合いを有している。したがって、発話 $U_j$ の時点でも、依然として「京都」は文脈上の意味合いに対して重みを持っている。なお、発話 $U_j$ を発したユーザAは少なくとも、祭に相当する単語として「祇園祭」などを想起しているはずである。

30

40

【0076】

これに対し、文書集合100中の文書101には京都の旅行記が記されている。その中の文単位 $S_i$ は、「7月」の「京都」といえば「祇園祭」という意味合いを有している。即ち、文単位 $S_i$ は、『「夏」の「7月」の「京都」の「祭」といえば』、「祇園祭」であるという意味合いを有している。つまり、発話 $U_j$ と、文単位 $S_i$ とは、共通して「夏」「京都」「祭」に重みを有しており、文脈上の意味合いが類似している。このように、本発明に係る文単位検索方法では、発話 $U_j$ の際にユーザが意識している、先行の発話か

50

らの文脈上の意味のまとまりを推定し、類似する文脈上の意味合いを有する文単位  $S_k$  を直接的に検索して出力することを目的としている。

【 0 0 7 7 】

本発明に係る文単位検索方法を実施するコンピュータシステムを実現した場合、連続する発話を受け付け、それらの言葉の文脈上の意味と類似する文単位を文書集合から抽出するのみならず、ユーザ A とユーザ B との会話中に、コンピュータシステムが発話毎に関連する情報を提示して会話に参入する鼎談が可能になる。また、コンピュータシステムがユーザ A とユーザ B との会話を支援することも可能になる。図 1 の説明図の例で、会話 100 のユーザ A による発話  $U_j$  の次に、コンピュータシステムによって「7月の京都といえば祇園祭です。」等の音声の出力がされた場合は、ユーザ A とユーザ B とコンピュータシステムとの間での鼎談が実現することになる。また、ユーザ A とユーザ B との会話が続かなくなった場合に、コンピュータシステムによって「7月の京都といえば祇園祭」等の情報の提示がされることで、ユーザ A とユーザ B との会話への支援も実現する。

10

【 0 0 7 8 】

そこで、このような文脈上の意味が類似する文単位を文書集合から検索することを実現するために、本発明に係る文単位検索方法をコンピュータ装置に実施させる。この場合、コンピュータ装置には、予め文書集合の文書データを夫々文単位に分別しておく処理、及び分別した文単位に各文単位の文脈上の意味を表わす定量的な情報を記憶させておく処理を含む事前処理が必要になる。さらに、コンピュータ装置が発話を受け付けた場合、その発話の会話の流れ上の意味を表わす定量的な情報を求める処理、及び、発話に対して求めた情報に基づいて意味が類似する文単位を抽出して検索結果として出力する処理を含む検索処理が必要になる。

20

【 0 0 7 9 】

したがって、以下に説明する実施の形態 1 乃至 3 では、本発明に係る文単位検索方法をコンピュータ装置に実施させるために必要なハードウェア構成についてまず説明する。さらにコンピュータ装置による処理を、事前処理と検索処理とを区別して段階的に説明する。具体的には、各実施の形態において、

「 1 . ハードウェアの構成及びシステムの概要」、

事前処理として

「 2 . 文書データの取得及び自然言語解析」、及び

「 3 . 文書データの文毎の意味のまとまりの定量化」、

次に

「 4 . 検索処理」

の順に説明する。

30

【 0 0 8 0 】

なお、以下に説明する実施の形態 1 乃至 3 では、本発明に係る文単位検索方法を実施する例として、文書データの文書集合を記憶しておくハードウェアと、発話を受け付けるコンピュータ装置と、文書集合が記憶されたハードウェア及び発話を受け付けるコンピュータ装置に接続して検索処理を実行するコンピュータ装置とで構成される検索システムを挙げて説明する。

40

【 0 0 8 1 】

また、以下に示す例では主に、文書集合が日本語の自然文からなる場合について各処理、具体例を示している。しかしながら、本発明の文単位検索方法は、日本語のみならず、他の言語にも適用することができることは勿論である。この場合、言語解析（形態素解析、統語解析）等の言語毎に特有の文法上の取り扱い等は、その言語毎に最適な方法を用いる。

【 0 0 8 2 】

（実施の形態 1）

1 . ハードウェアの構成及びシステムの概要

図 2 は、実施の形態 1 における文単位検索装置 1 を用いた検索システムの構成を示すブ

50

ロック図である。検索システムは、文書データからの検索処理を実行する文単位検索装置 1 と、自然言語からなる文書データを記憶する文書記憶手段 2 と、インターネット等のパケット交換網 3 と、ユーザから入力されるキーワード又は音声等の言葉を受け付ける受付装置 4, 4, ... とで構成される。文単位検索装置 1 は、PC (Personal Computer) であり、自然言語からなる文書データを記憶する文書記憶手段 2 と接続される。また、受付装置 4, 4, ... も PC であり、文単位検索装置 1 は、パケット交換網 3 を介して受付装置 4, 4, ... と接続され通信が可能である。

【0083】

実施の形態 1 の検索システムでは、文単位検索装置 1 は、検索の対象である文単位を含む文書データを文書記憶手段 2 に予め記憶しておく。文単位検索装置 1 は、文書記憶手段 2 に記憶した文書データを、予め文単位に分別し、後に検索処理が可能ないように各文単位に文脈上の意味を表わす定量的な情報を記憶させておく。また、受付装置 4, 4, ... は、受け付けた言葉をコンピュータで処理可能なテキストデータ又は音声データに変換し、パケット交換網 3 を介して当該データを文単位検索装置 1 へ送信する。文単位検索装置 1 が、受信した言葉のデータに基づいて文書記憶手段 2 に記憶した文書データから一又は複数の文からなる文単位を抽出し、抽出した文単位をパケット交換網 3 を介して受付装置 4, 4, ... へ出力することで文単位の検索を実現する。

【0084】

文単位検索装置 1 は、少なくとも、各種ハードウェアを制御する CPU 11 と、各種ハードウェア間を接続する内部バス 12 と、不揮発性のメモリからなる記憶手段 13 と、揮発性のメモリからなる一時記憶領域 14 と、パケット交換網 3 と接続するための通信手段 15 と、文書記憶手段 2 と接続するための文書集合接続手段 16 と、DVD、CD-ROM 等の可搬型記録媒体 18 を用いる補助記憶手段 17 とを備える。

【0085】

記憶手段 13 には、DVD、CD-ROM 等の可搬型記録媒体 18 から取得した、PC が本発明に係る文単位検索装置 1 として動作するための制御プログラム 1P が記憶されている。CPU 11 は、制御プログラム 1P を記憶手段 13 から読み出して実行すると共に、内部バス 12 を介して各種ハードウェアを制御する。一時記憶領域 14 は、CPU 11 の演算処理によって一時的に発生する情報が記憶される。

【0086】

CPU 11 は、受付装置 4, 4, ... から送信される言葉のデータを通信手段 15 を介して受信したことを検知し、受信した言葉のデータに基づいて処理を実行し、検索処理を行う。また、CPU 11 は、文書集合接続手段 16 を介して文書記憶手段 2 で記憶している文書データを取得し、且つ、文書集合接続手段 16 を介して文書データを文書記憶手段 2 に記憶させることが可能である。

【0087】

DVD、CD-ROM 等の可搬型記録媒体 18 から補助記憶手段 17 を介して取得した、記憶手段 13 に記憶されている制御プログラム 1P では更に、記憶手段 13 で記憶している辞書情報に基づいて文字列で表された文書データを形態素解析及び統語解析等の自然言語解析を CPU 11 に実行させることができるようにしてある。

【0088】

受付装置 4, 4, ... は、少なくとも、各種ハードウェアを制御する CPU 41 と、各種ハードウェア間を接続する内部バス 42 と、不揮発性メモリからなる記憶手段 43 と、揮発性メモリからなる一時記憶領域 44 と、マウス又はキーボード等の操作手段 45 と、モニタ等の表示手段 46 と、マイク及びスピーカ等の音声入出力手段 47 と、パケット交換網 3 へ接続するための通信手段 48 とを備える。

【0089】

記憶手段 43 には、PC が受付装置 4, 4, ... として動作するための処理プログラム等が記憶されている。CPU 41 は、処理プログラムを記憶手段 43 から読み出して実行すると共に、内部バス 42 を介して各種ハードウェアを制御する。一時記憶領域 44 は、C

10

20

30

40

50

P U 4 1 の演算処理によって一時的に発生する情報が記憶される。

【 0 0 9 0 】

C P U 4 1 は、ユーザからの文字列入力操作を操作手段 4 5 を介して検知し、入力された文字列を一時記憶領域 4 4 に記憶することができる。C P U 4 1 は、ユーザから入力された音声を音声入出力手段 4 7 を介して検知し、記憶手段 4 3 に記憶された音声認識のためのプログラムを読み出して実行することによって入力された音声をテキストデータに変換することができる。また、C P U 4 1 は、ユーザから入力された音声を音声入出力手段 4 7 により、コンピュータで処理可能な音声データとして入力することができる。

【 0 0 9 1 】

また、C P U 4 1 は、ユーザからの文字列入力操作又は音声入力を検知することで得られたテキスト又は音声の言葉のデータを通信手段 4 8 を介して文単位検索装置 1 へ送信する。

10

【 0 0 9 2 】

なお、C P U 4 1 は、音声データをテキストデータに変換して送信してもよく、その場合は、C P U 4 1 は、音声認識によって得られる音声データの特徴、例えば各単語に相当する音素が発声された時の速度、単語に相当する音素の周波数等のデータを共に送信してもよい。また、C P U 4 1 は、各単語に相当する音声データ間の時間差についても記憶しておき、以前に受け付けた言葉にその単語が含まれていた時点との時間差も共に文単位検索装置 1 へ送信してもよい。

【 0 0 9 3 】

20

## 2 . 文書データの取得及び自然言語解析

上述のように構成される検索システムにおいて、文単位検索装置 1 はまず、事前処理として文書集合を用意して、後に各文書データに含まれる文単位毎の意味のまとまりを表わすことができるようにしておく処理を行なう。「2 . 文書データの取得及び自然言語解析」では、文単位検索装置 1 が文書記憶手段 2 に文書データを記憶しておき、各文書データを言語解析して一又は複数の文からなる文単位に分別し、さらに文単位毎に文法的な特徴を解析し、文書記憶手段 2 に文単位毎に記憶しておく処理について説明する。なお、実施の形態 1 では、文単位検索装置 1 は文単位を一の文とした場合について説明する。

【 0 0 9 4 】

文単位検索装置 1 の C P U 1 1 は、検索の対象である文単位を含む文書データを文書記憶手段 2 に予め記憶しておく。文単位検索装置 1 の C P U 1 1 は、通信手段 1 5 及びパケット交換網 3 を介して取得可能な文書データを Web クローリングにより取得し、文書集合接続手段 1 6 を介して文書記憶手段 2 に記憶する。文単位検索装置 1 の C P U 1 1 は、取得して文書集合接続手段 1 6 を介して文書記憶手段 2 に記憶してある文書データを文単位に分別し、夫々言語解析（形態素解析及び統語解析）を行い、その結果を文単位毎に対応付けて記憶する処理を行なう。

30

【 0 0 9 5 】

以下に、文単位検索装置 1 の C P U 1 1 が、文書データを取得し、取得した文書データに対して形態素解析及び統語解析の自然言語解析をして、文単位毎に記憶する処理手順について説明する。図 3 は、実施の形態 1 における文単位検索装置 1 の C P U 1 1 が、取得した文書データに対する形態素解析及び統語解析処理の解析結果からタグ付け及び単語抽出を行い記憶する処理手順を示すフローチャートである。図 3 のフローチャートに示す処理は、文単位毎にその文単位に出現する単語又は先行の文単位から参照する単語を抽出する処理と、各文単位における各単語の特徴を特定して記憶しておく処理に対応する。

40

【 0 0 9 6 】

C P U 1 1 は、Web クローリングを開始すると文書データを取得したか否か判断する（ステップ S 1 1）。C P U 1 1 が文書データを取得していないと判断した場合は（S 1 1 : N O）、C P U 1 1 は処理をステップ S 1 1 へ戻し、文書データを取得するまで待機する。C P U 1 1 が文書データを取得したと判断した場合は（S 1 1 : Y E S）、C P U 1 1 は、取得した文書データから一文毎の読み出しを試み、読み出しが成功したか否かを

50

判断する(ステップS12)。

【0097】

CPU11が、読み出し箇所が文書データの終端に至っておらず、文の読み出しが成功したと判断した場合は(S12: YES)、読み出した文の形態素解析及び統語解析を行う(ステップS13)。

【0098】

CPU11は、形態素解析及び統語解析の結果から、解析した文に出現する単語及び当該文で先行の文から参照する単語を抽出し、リストに記憶する(ステップS14)。更に、CPU11は、後述で説明するように解析結果からタグを生成し(ステップS15)、読み出した文にタグを付加して、文書集合接続手段16を介して文書記憶手段2に記憶させる(ステップS16)。

10

【0099】

一方、CPU11が、読み出し箇所が文書データの終端に至っており、文の読み出しが失敗したと判断した場合は(S12: NO)、取得した文書データに対する処理を終了する。

【0100】

上述の処理を、文書データを取得する都度に行い、タグ付け済みの文書データを文書記憶手段2に記憶しておく。

【0101】

次に、文単位検索装置1のCPU11による上述の処理の詳細を、具体例を挙げて説明する。

20

【0102】

図4は、実施の形態1における文書記憶手段2で記憶される文書データの内容の一例を示す説明図である。文書記憶手段2で記憶される文書データは、文単位検索装置1のCPU11が通信手段15を介して、パケット交換網3に接続され公開されているWebサーバから取得されたHTML(HyperText Markup Language)等のテキストデータをもとに記憶される。図4に示す一例も、インターネットで公開されたWebページ(<http://ja.wikipedia.org/wiki/祭より抜粋>)より取得することができたHTMLデータの文書である。以下、この文書例を使用して文書の解析及び検索等について説明する。

30

【0103】

文単位検索装置1のCPU11は、図3のフローチャートに示したステップS12の文の読み出しの処理において、取得した文書データ中の文字列を「文」の言語単位(文単位)に分別する。分別する方法として例えば、CPU11は、日本語からなる文書データである場合、句点「。」を表す文字列によって、又は、英語からなる文書データである場合はピリオド「.」を表す文字列によって分別してもよい。

【0104】

次に、図3のフローチャートに示した文単位検索装置1のCPU11によるステップS13の形態素解析及び統語解析の処理の詳細を説明する。

【0105】

文単位検索装置1のCPU11は、「文」の言語単位に対して辞書情報に基づいた形態素解析を行い、文の最小構成単位である形態素を同定して形態素の構造を解析する。例えば、図4に示した文書データでは、CPU11は、記憶手段13の辞書情報に基づいて、「祭」「神霊」等の名詞、「九州」等の固有名詞、「祀る」等の動詞、「と」「は」等の助詞、「、」「。」等の記号等を示す文字列と照合することで形態素を同定する。形態素解析の手法については今日では種々の手法が提案されており、本発明では当該形態素解析の手法を限定するものではない。

40

【0106】

さらに、文単位検索装置1のCPU11は、同定した形態素毎にその品詞情報(名詞、助詞、形容詞、動詞、副詞等)と、日本語文である場合は日本語の文法、英文である場合

50

は英語の文法に基づく品詞間の結束性を統計的に求めた文法情報とに基づいて形態素間の文法的関係を抽出する統語解析を行う。例えば、文法を木構造に当てはめて形態素の品詞情報から木構造に従って形態素間の関係を抽出することができる。解析対象が(形容詞+名詞+助詞+名詞)である場合、まず解析対象が名詞であるか否かを判断する。名詞でないと判断した場合は次に、当該解析対象が(形容詞+名詞)に当てはまるか否かを判断する。したがって、当該解析対象の先頭の形態素が形容詞句であるか否かを判断する。先頭の形態素が形容詞であると判断した場合は、当該形容詞が後続する名詞を修飾する当該解析対象の中で一番大きな修飾語であると判断される。つまり(形容詞+(名詞))という関係が抽出される。

【0107】

次に、残りの解析対象が(名詞)であるか否かを判断する。複数の形態素からなり、名詞ではないと判断した場合は、当該残りの解析対象が(形容詞+名詞)に当てはまるか否かを判断する。したがって、残りの解析対象の先頭の形態素が形容詞であるか否かを判断する。残りの解析対象の先頭の形態素が形容詞でないと判断した場合は、(形容詞+名詞)の形容詞の部分(名詞+助詞)に展開し、残りの解析対象が((名詞+助詞)+名詞)に当てはまるか否かを判断する。残りの解析対象が((名詞+助詞)+名詞)に当てはまると判断した場合は、当該解析対象(形容詞+名詞+助詞+名詞)の形態素間の文法的関係は[形容詞+{(名詞+助詞)+名詞}]であると抽出することができる。統語解析の方法についてもこのような方法を基礎とする手法に限らず、形態素解析の手法同様に今日では種々の手法が提案されており本発明では当該統語解析の手法を限定するものではない。

【0108】

実施の形態1では、一例として形態素解析及び統語解析についてchasen(<http://chasen.org>)及びCabocha(工藤 拓、松本 裕治「チャンキングの段階適用による日本語係り受け解析」情報処理学会論文誌Vol.6、No.43、pp.1834-1842(2002)、<http://chasen.org/~taku/software/cabocha>参照)にて開示された技術に基づいて行う。他にKNP(Kurohashi-Nagao Parser)(黒橋 禎夫、長尾 眞「並列構造の検出に基づく長い日本語文の構造解析」自然言語処理Vol.1、No.1、pp.35-57(1994))で開示されている技術に基づいて解析するのもよい。

【0109】

文単位検索装置1のCPU11は、解析した形態素及び形態素間の文法的関係を、XML(Extensible Markup Language)に基づくタグで表した文書データを生成して文書記憶手段2に記憶させる。本発明が利用する形態素解析及び統語解析の自然言語解析方法(chasen、Cabocha)では入力された文字列を形態素解析し、さらに統語解析して各形態素の品詞情報、形態素の係り先を示す情報等を分別した形態素毎に出力するようにしてある。文単位検索装置1の記憶手段13に記憶されている制御プログラム1Pでは、当該自然言語解析方法を文単位検索装置1のCPU11に実行させることができるように構成されている。

【0110】

本発明が利用する形態素解析及び統語解析では、例えば、図4に示した「九州地方北部では、秋に行われるものに対して(お)くんちと称する場合もある。」という文の文字列に対しまず文節番号が付される。(0:九州地方北部では、/1:秋に行われるものに対して(お)くんちと称する場合も/2:ある。)さらに各文節で形態素に分別され、形態素毎の品詞情報、形態素の基本形情報、発音情報等が付加される。文節番号0の文節は、(0:九州(名詞+固有名詞+地域+一般、九州、キュウシュウ)/地方(名詞+一般、地方、チホウ)/北部(名詞+一般、北部、ホクブ)/で(助詞+格助詞+一般、で、デ)/は(助詞+係助詞、は、ハ)/、(記号+読点))と形態素の同定及び情報の付加が行われる。「九州」という形態素は名詞であって固有名詞であり、地域を示す名詞でもあり、一般名詞として使用されることもある。また基本形は「九州」であり、「キュウシュ

10

20

30

40

50



ウ」と発音することを判別することができる。他の文節も同様である。また、係り受け情報は例えば、(0 2, 1 2, 2 - 1)と文節間の係り受け関係が判別可能なように取得できる。この例では、文節番号0の文節は文節番号2の文節を係り先とし、文節番号1の文節は文節番号2の文節を係り先とすることが判別できる。また、文節番号2の文節は係り先がないことを係り先を-1とすることで判別できる。

【0111】

図5は、実施の形態1における文単位検索装置1のCPU11が、形態素解析及び統語解析した結果を付与して文書記憶手段2に記憶させる文書データの一例を示す説明図である。図4に示した内容の文書データに対して図3のフローチャートに示した処理手順が実行されたことにより文書記憶手段2に記憶された文書データの例に相当する。

10

【0112】

図5に示すように、文単位検索装置1のCPU11により、図4に示した内容の文書の一部が固有名詞、名詞、助詞、動詞等の形態素に分別され、形態素間の文法的関係性はタグの入れ子によって表されている。図5に示す例は、GDA(Global Document Annotation; <http://i-content.org/gda>参照)で提案されている規則に則ったタグ付け手法に従ったものである。本発明では当該規則に従うことを限定するものではない。また、形態素の情報及び形態素間の係り受けの情報をコンピュータが情報処理によって識別できるようにすることができればXMLのタグ付けによる方法には限らない。

【0113】

GDAに基づくタグ付けは基本的に<タグ名 属性名 = “属性値”>で表される。図5に示される例では、<su>で示されるタグは、文(Sentential unit)を表すタグである。図5に示した例では、「九州地方北部では、秋に行われるものに対して(お)くんちと称する場合もある。」の文は、「九州地方北部では」「、」「秋に行われるものに対して(お)くんちと称する場合もある」「ある」「。」の三つの文節と句読点との単位を有していることがタグによって判別できる。<ad>で示されるタグは、終助詞以外の助詞(particle)、副詞(adverb)、連体詞などを示すタグであるが、文節0の「九州地方北部では」も全体で副詞的な役割を果たすことを示すことができる。<n>で示されるタグは、名詞(noun)を示す。<v>で示されるタグは、動詞(verb)を示す。また、図5に示したタグの他に形容詞(adjective)を示す<aj>タグ等がある。

20

30

【0114】

属性名synで表される属性は、当該属性が付与されているタグで挟まれた文節又は語等の言語単位間の係り受け関係を示す。属性値f(forward; 前向き)が付与されている文では、当該文を構成する言語単位は一番近い後続の言語単位に係ることを示す。したがって、原則では文節0の「九州地方北部では」は、文節1の「秋に行われるものに対して(お)くんちと称する場合もある」へ係り、文節1の「秋に行われるものに対して(お)くんちと称する場合もある」は文節2の「ある」に係る。

【0115】

しかし統語解析により、文節0の「九州地方北部では」は文節2の「ある」に係り、文節1の「秋に行われるものに対して(お)くんちと称する場合もある」は文節2の「ある」に係ることが判別できているため、上述原則はあてはまらない。したがって、係り受けの受ける側ではない「句」(phrase)であることを示す“p”を各タグに付加することで、係り受けの関係を示すことができる。例えば、<adp>で示されるタグは、タグ<ad>に、句であることを示す“p”が組み合わさったものである。<adp>タグでは含まれた文節は副詞句であって、係り受けの受ける側の文節ではないことを示す。したがって、図5に示した例では、文節1の「秋に行われるものに対して(お)くんちと称する場合もある」は、副詞句であって受ける側の文節ではないため、文節0の「九州地方北部では」は、文節1の「秋に行われるものに対して(お)くんちと称する場合もある」へ係らずに「ある」に係ることが示される。その他、“p”は「句」であることを明示するために付加

40

50

される。

【0116】

また、< n >で示すタグについても、< n p >とすることで係り受けの受ける側の語ではないことを示すことができる。「九州地方北部」は、「九州」「地方」「北部」と夫々< n >で挟まれる形態素に分別でき、「九州」は「地方」に、「地方」は「北部」に係るため“ p ”は不要である。一方、「催事（催し、イベント）、フェスティバルのこと」では、「催事（催し、イベント）」は「フェスティバル」に係らず「の」に係るため、「フェスティバル」を挟むタグを< n p >とすることで、係り受けの関係を示すことができる。

【0117】

なお、「九州」のような場所を表す固有名詞、又は「太郎」のような人の名前を表す固有名詞は、夫々< p l a c e n a m e >< p e r n a m e >のタグによって示すことができる。

【0118】

指示代名詞、ゼロ代名詞等の先行する語又は文から参照する形態素については、照応関係を表す属性を用いて表すことができる。GDAでは、属性名idを用いて指示代名詞、ゼロ代名詞が先行の語又は文の何れの語を示すかをあらわすことができる。例えば、「右側にボタンがあるので、それを押してください。」という文に対して、人間がこれを読む場合は「それ」が「ボタン」を指すことを自然に補完することができる。しかし、コンピュータで処理する場合は、辞書情報との照合によって「それ」が指示代名詞であることを同定することはできるが、何を示しているかを判別することはできない。そこでGDAでは、「それ」が示す「ボタン」にid属性を付加し、さらに、id属性で示された形態素との等価(e q u a l)関係を示す属性名eqにより、「それ」=「ボタン」を示すことができる。具体的には「右側にボタンがあるので、それを押してください。」に対し、「右側に< n p i d = “ B t n ” >ボタン< / n p >があるので、< n p e q = “ B t n ” >それ< / n p >を押してください。」とすることで(他のタグは省略)、「それ」=「ボタン」の関係を示すことができる。

【0119】

ゼロ代名詞に対しては、eq属性を付加できる代名詞そのものがない。したがって、「それ」=「ボタン」を動作の対象とする「押し」という動詞に、対象を明示する情報を付加することで、ゼロ代名詞が表す対象を示すことができる。そこで、タグではさんだ形態素の動作の対象(o b j e c t)を示す属性名objにより、「押し」という動作の対象が「ボタン」であることを示すことができる。具体的には、「右側にボタンがあるので、押してください。」という文に対し、「右側に< n p i d = “ B t n ” >ボタン< / n p >があるので、< v o b j = “ B t n ” >押し< / v >てください。」とすることで、省略された対象との関係を明示することができる。

【0120】

また、参照される語と参照する語とが離れている場合であっても、上述のid属性、eq属性、obj属性によってその照応関係を示すことができる。例えば、「右側に< n p i d = “ B t n ” >ボタン< / n p >があります。」「< n p e q = “ B t n ” >それ< / n p >にはxのマークがついています。」「停止する際に< v o b j = “ B t n ” >押し< / v >てください。」とすることによって、第2文の「それ」が「ボタン」を示すこと、及び第3文の「押し」の対象が「ボタン」であることを示すことができる。

【0121】

また、各形態素を挟む< n >< a d >< v >等のタグの属性情報には、形態素(m o r p h e m e)解析の結果を示す情報が属性名m p hで付加される。属性値は、形態素解析によって取得できた形態素の品詞情報、基本形情報、発音情報等を示す。具体的には、属性名m p hに対し、付加情報、品詞情報、活用形情報、基本形情報、及び発音情報を属性値とし、m p h = “付加情報；品詞情報；活用形情報；基本形情報；発音情報”と表す。図5に示した例において「九州」は、品詞情報を名詞+固有名詞+地域+一般で分類す

10

20

30

40

50

ることができ、基本形は九州であり「キュウシュウ」と発音することが < m p h > タグによって明示される。なお、本発明では、形態素解析及び統語解析を c h a s e n で提示される方法に基づいて行っているため、形態素の付加情報として c h a s e n という識別情報が付加されている。

【 0 1 2 2 】

上述のように、文単位検索装置 1 の C P U 1 1 は W e b クローリングによって取得した文書データに対し、形態素解析及び統語解析の結果を G D A の規則に則ってタグ付けし、タグ付けした結果である X M L データを文書集合接続手段 1 6 を介して文書記憶手段 2 に記憶させる。文書データを X M L データで記憶しておくことにより、文単位検索装置 1 の C P U 1 1 は当該文書データのタグを文字列解析によって識別し、タグに付加された属性情報を識別することによって各形態素の情報及び文法的関係を特定することができる。

10

【 0 1 2 3 】

さらに文単位検索装置 1 の C P U 1 1 は、W e b クローリングによって取得した文書データを形態素解析する際に、取得した全文書データに出現する単語を抽出して識別番号を割り振りリストで記憶手段 1 3 に記憶する。図 6 は、実施の形態 1 における文単位検索装置 1 の C P U 1 1 が取得した全文書データから抽出した単語のリストの例を示す説明図である。図 6 の説明図に示す例では、3 1 2 4 5 個の単語がリストとして挙げられている。なお、記憶される単語からは、「こと」、「もの」などのありふれた単語は除かれる。接続詞又は冠詞同様一般的すぎる言葉であり、頻繁に出現するにも拘わらず、その単語自体は意味をなさないために検索処理に負担がかかり、検索対象として不適切であるからである。

20

【 0 1 2 4 】

3 . 文書データの文毎の意味のまとまりの定量化

3 - 1 . 文毎の意味のまとまりの定義

次に、文単位検索装置 1 の C P U 1 1 は、文書記憶手段 2 で記憶した文書データ中の一文毎に当該文の意味のまとまりを定量的に表す情報を特定する。文の意味のまとまりを定量的に表す情報とは、ユーザが当該文を使用（発話、筆記、聴取又は読解）するときに、ユーザが注目している単語群と、ユーザが各単語に注目する度合い、即ち顕現性（s a l i e n c e）を定量的に示す値（単語の重み値）とで表す。

【 0 1 2 5 】

各単語の文中での顕現性は、従来の検索サービスによってされてきた出現頻度によって定量化することもできる。しかしながら、出現頻度は文書、又は文書集合全体を母体として求めるものである。したがって、文書毎に各単語の出現頻度を算出することで、文書全体の意味のまとまりを定量的に表すことはできても、文書中での流れに応じて一文毎に動的に変化する文脈を反映した意味のまとまりを表すことはできない。

30

【 0 1 2 6 】

また、単語の文中での顕現性は、先行する文での当該単語の注目度、現在の文での当該単語の注目度の遷移をその単語の使用のされ方で文法的に区別して表すことができる。つまり、先行する文で主題（主語）であった単語が現在の文でも主題（主語）である場合は、現在の文で当該単語は一番注目されている顕現性の高い単語であるといえる。これに対し先行する文では出現していないが現在の文で主題（主語）である単語は、現在の文で注目されているものの、前述の主題として使用され続ける場合に比べて顕現性は低いといえる。この顕現性の定式化は、中心化理論（Grosz et al., 1995、Nariyama, 2002、Poesio et al., 2004）として研究が続けられている。

40

【 0 1 2 7 】

中心化理論による定式化では、各単語の顕現性をコンピュータ等で定量的に計算するための特徴量として表わされていない。各単語の遷移の仕方が中心化理論で定義される遷移の仕方の何れに属するか否かが判別できるに過ぎない。そこで本発明では各単語の各文での顕現性を定量的に算出する。

【 0 1 2 8 】

50

実施の形態1では、単語毎に各文単位での参照確率を算出し、算出した参照確率を各単語の文単位での顕現性を表わす重み値として付与する。

【0129】

なぜなら、単語が当該文で注目されているほど、継続して後続の文でも出現又は参照される確率が高いことから、後続の文で出現する確率又は後続の文から参照される確率を参照確率とし、当該単語の顕現性と捉えることができるからである。また、単語が後続の文で出現又は参照される参照確率は、定量的に扱うことが困難な単語の意味を特徴とするのではなく、文単位検索装置1による情報処理によって解析可能な、単語が出現するパターン又は参照するパターンを含む特徴パターンを特定し、特定した特徴パターンと同一の特徴パターンで出現又は参照される単語が実際に後続の文で出現又は参照される割合が参照確率として算出される。

10

【0130】

以下、単語毎の参照確率を各単語の文単位での重み値とし、夫々の重み値が付与された当該文での単語の集合を重み付き単語群という。各文単位の意味のまとまりは、参照確率という定量的な重み値が付与された重み付き単語群で表わすことができる。

【0131】

### 3-2. 回帰モデル学習

参照確率の算出は、特定した特徴パターンと同一の特徴パターンが出現した数に対して、同一の特徴パターンのうち当該単語が実際に後続の文で出現又は参照される割合をその参照確率として求める。この際、特定した特徴パターンと同一の特徴パターンが夫々の特徴パターン毎に多量に且つほぼ同数で出現する場合は、統計的に問題なく参照確率を算出することができる。しかし、実際に同一の特徴パターンが出現する数は限られ、信頼に足り得る参照確率を算出するには膨大な文書データが必要となる。したがって、後続の文で出現又は参照されるか否かをその事象の発生の要因である単語の特徴パターンから予測するための回帰式を、特徴パターンと実際に後続の文で出現又は参照されたかの事象とで回帰モデル学習をすることによって求める。

20

【0132】

以下、回帰モデル学習のためのサンプルである特徴パターンに対する「3-2-1. 特徴パターンの特定」と特徴パターンを用いた「3-2-2. 回帰式の学習」とに段階を分けて説明する。

30

【0133】

### 3-2-1. 特徴パターンの特定

文書記憶手段2に記憶してある文書データ中の文は< s u >で示すタグで挟まれ、当該文で出現する単語、若しくは文の中の指示代名詞又はゼロ代名詞と照応関係にある単語は、タグの属性情報により特定することが可能である。そこで、本発明の文単位検索装置1では、文書記憶手段2で記憶した文書データに対し、特徴パターンを以下のように特定する。

【0134】

文書データ中の一の文sと、当該文書データ中での一の文に対する先行する文に含まれる単語wの対をサンプル(s, w)とする。当該サンプルに対する特徴パターンf(s, w)は、以下の特徴量によって特定される。文sと、文sより先行する文のうち単語wが、最近に出現又は参照された文sとの距離(文の数)の特徴量(dist)、文sより先行する文で単語wが、最近に出現又は参照された場合、単語wが係っている助詞の特徴量(gram)、及び文sより先行する文で単語wが出現又は参照された数(chain)の特徴量等を例として挙げるることができる。なお、特徴量はこれに限らず、単語wが最近のトピックを示す単語であるか否か、又は単語wが一人称であるか否か等でもよい。

40

【0135】

文書記憶手段2で記憶した文書データには形態素解析及び統語解析の結果がGDAに則ったタグによって記述されているため、文書データの文字列解析によってタグ< s u >で区切られる文の分別及び計数、各文内のタグで示される品詞情報による助詞の特定、指示

50

代名詞又はゼロ代名詞で参照するものも含んだ単語の出現回数の計数が可能である。したがって、文単位検索装置 1 の CPU 1 1 は、GDA に則ったタグ及びその属性値を解析することで各サンプルに対する特徴量  $d i s t , g r a m , c h a i n$  を特定することができる。

#### 【0136】

文単位検索装置 1 の CPU 1 1 が、文書記憶手段 2 で記憶しているタグ付け済みの文書データに対しサンプルを抽出し、抽出したサンプルに対して特徴量を求めて特徴パターンを特定し、抽出したサンプルの特徴パターンから参照確率を算出するため回帰式を回帰分析により推定する処理手順について説明する。図 7 は、実施の形態 1 における文単位検索装置 1 の CPU 1 1 が、文書記憶手段 2 で記憶しているタグ付け済み文書データからサンプルを抽出し、回帰分析を行って参照確率を算出するための回帰式を推定する処理手順を示すフローチャートである。図 7 のフローチャートに示す処理は、分別した文単位毎に特徴パターンを特定する処理、及び、特徴パターンと、特定された単語が後続の文単位で出現又は参照されたか否かの判定結果とに基づいて参照確率を算出するための回帰学習を実行する処理に対応する。

10

#### 【0137】

文単位検索装置 1 の CPU 1 1 は、文書記憶手段 2 から文書集合接続手段 1 6 を介してタグ付け済みの文書データを取得する (ステップ S 2 1 )。CPU 1 1 は、取得した文書データに付加されたタグ  $\langle s u \rangle$  を文字列解析によって識別して文に分別する (ステップ S 2 2 )。次に CPU 1 1 は、文を示す  $\langle s u \rangle$  内の各タグを文字列解析によって識別し、文に対し当該文で出現する単語又は参照される単語を対応付けてサンプルを抽出する (ステップ S 2 3 )。抽出したサンプルに対し、タグを文字列解析によって識別して  $d i s t , g r a m , c h a i n$  からなる特徴パターンを特定する (ステップ S 2 4 )。

20

#### 【0138】

CPU 1 1 は、分別した文が取得した文書データの終端であるか否かを判断し (ステップ S 2 5 )、CPU 1 1 が、分別した文が文書データの終端でないと判断した場合は (S 2 5 : N O)、CPU 1 1 は処理をステップ S 2 2 に戻し、後続の文について  $\langle s u \rangle$  タグを識別することで分別する処理を継続する。分別した文が取得した文書データの終端であるか否かは、例えば現在分別した文を挟む  $\langle s u \rangle \langle / s u \rangle$  の後に、 $\langle s u \rangle$  タグが後続するかないかを判断し、後続しないと判断した場合は終端であると判断することができる。

30

#### 【0139】

一方、CPU 1 1 が文書データの終端であると判断した場合は (S 2 5 : Y E S)、CPU 1 1 は、所定の数のサンプルの抽出が終了したか否かを判断する (ステップ S 2 6 )。CPU 1 1 がサンプルの抽出が終了していないと判断した場合は (S 2 6 : N O)、CPU 1 1 は、処理をステップ S 2 1 へ戻し、異なるタグ付け済みの文書データを取得し、サンプルの抽出を継続する。

#### 【0140】

CPU 1 1 がサンプルの抽出が終了したと判断した場合は (S 2 6 : Y E S)、CPU 1 1 は、抽出したサンプルに対して回帰分析を行い、各特徴量  $d i s t , g r a m , c h a i n$  に対する回帰式の回帰係数を推定し (ステップ S 2 7 )、処理を終了する。

40

#### 【0141】

次に、文単位検索装置 1 の CPU 1 1 による上述の処理の詳細を、具体例を挙げて説明する。

#### 【0142】

図 8 は、実施の形態 1 における文書記憶手段 2 で記憶された文書データ中の文で特定される特徴パターンの例を示す説明図である。図 8 に示す文  $s_i$  での、当該文  $s_i$  と、先行する文に含まれる単語「太郎君」とのサンプル ( $s_i$ , 太郎君) の特徴パターン  $f(s_i, \text{太郎君})$  は以下のようにして特定される。現在の文  $s_i$  と、先行する文のうち最近に、単語「太郎君」が出現又は参照された文  $s_{i-1}$  との距離の特徴量 ( $d i s t$ ) は、 $s_i$  の

50

直後に続く文  $s_{i+1}$  までの文の数 2 であるため  $dist = 2$  である。また、最近「太郎君」が出現又は参照された  $s_{i-1}$  での単語「太郎君」（彼で参照）が係っている助詞は「は」であるため、 $gram = 八$  である。更に、文  $s_i$  より先行の文  $s_{i-2}$ 、 $s_{i-1}$  で単語「太郎君」が出現又は参照されたため  $chain = 2$  である。したがって、特徴パターンは  $f(s_i, \text{太郎君}) = (dist = 2, gram = 八, chain = 2)$  と特定される。英語の場合、 $gram$  は前置詞によって特定される。

#### 【0143】

上述のように、文書データ中の文からサンプル  $(s, w)$  を抽出し、抽出した全サンプルに対して特徴パターン  $f(s, w)$  を特定する。

#### 【0144】

##### 3-2-2. 回帰式の学習

次に、図 7 のフローチャートに示したステップ S27 の回帰分析について、詳細な処理を説明する。

#### 【0145】

実施の形態 1 では、Logistic Regression モデルに基づいて回帰分析を行う。回帰分析はこれに限らず、kNN (k-Nearest Neighbors) 平滑化 + Support Vector Regression (SVR) モデルなど、他の回帰分析の手法を使用してもよい。

#### 【0146】

kNN 平滑化 + SVR モデルを使用する場合、扱うことのできる特徴パターンの特徴量として、次の 8 要素を使用して回帰モデルの学習ができる。8 要素とは、前述の  $dist$ 、 $gram$ 、 $chain$  に加えて、以下の 5 要素を特徴量として扱うことができる。一つは、先行の文単位の中で単語  $w$  を参照した場合の名詞の種別 ( $exp$ , 代名詞: 1 / 非代名詞: 0) でもよい。また、他の一つは、その単語  $w$  が先行の文単位において出現又は参照されている場合に主題であるか否か ( $last\_topic$ ,  $yes: 1 / no: 0$ ) でもよい。他の一つは単語  $w$  が先行の文単位において出現又は参照されている場合に主語であるか否か ( $last\_subj$ ,  $yes: 1 / no: 0$ ) でもよい。他の一つは、サンプル  $(s, w)$  において、単語  $w$  が一人称であるか否か ( $p1$ ,  $yes: 1 / no: 0$ ) でもよい。他の一つは、単語  $w$  が出現又は参照されている直近の先行の文単位での単語  $w$  の品詞情報 ( $pos$ , 名詞: 1、動詞: 2、等) でもよい。さらに他の一つは、単語  $w$  が文書中のタイトル又は見出しで参照されているか否か ( $in\_header$ ,  $yes: 1 / no: 0$ ) でもよい。さらに、音声データに基づいて回帰分析する場合、8 要素の内の 1 つとして、単語の直近の参照箇所の発話時刻からの秒数 ( $time\_dist$ )、単語の直近の参照箇所を含む文節の 1 音節あたりの発話速度 (の話者平均に対する比) ( $syllable\_speed$ )、単語の直近の参照箇所を含む文節の、最低発話音高と最高発話音高の周波数比 ( $pitch\_fluct$ ) の内のいずれか一又は複数を使用することができる。音声データの特徴量についても回帰分析することにより、後述するように文単位検索装置 1 の CPU 11 が言葉のデータとして音声データを受信した場合に、その特徴量から参照確率を算出することができる。

#### 【0147】

このように、kNN 平滑化 + SVR モデルを使用する場合、より詳細な特徴量に基づいて参照確率を算出することができ、より緻密な参照確率を算出することができる。

#### 【0148】

本実施の形態 1 では、文  $s_i$  の後続の文  $s_{i+1}$  で単語  $w$  が実際に出現又は参照されたか否かを被説明変数、サンプル  $(s_i, w)$  に対して特定された特徴パターンの  $dist$ 、 $gram$ 、 $chain$  を特徴量とし、全サンプル  $(s, w)$  に対して、Logistic Regression モデルにより回帰分析する。これにより、 $dist$ 、 $gram$ 、 $chain$  という特徴量が与えられた場合に、 $s_{i+1}$  で単語  $w$  が出現又は参照される確率  $Pr(s_{i+1}, w)$  を算出するための回帰式を得ることができる。

#### 【0149】

10

20

30

40

50

Logistic Regressionモデルで求められる確率は、一般的に、説明変数（特徴量） $x_1, x_2, \dots, x_n$ に対して以下の式（1）で求められる。

【0150】

【数1】

$$\text{Pr} = \frac{1}{1 + \exp(b_0 + b_1x_1 + b_2x_2 + \dots + b_nx_n)} \quad \dots (1)$$

10

【0151】

式（1）のパラメータ（回帰係数） $b_0, b_1, \dots, b_n$ は、学習するサンプルから最尤法によって推定する。本発明で算出する文 $s$ での単語 $w$ の参照確率の回帰分析とは、被説明変数を、後続の文 $s_{i+1}$ で出現又は参照されないサンプルは0、出現又は参照されるサンプルは1とし、説明変数を特徴量である $dist, gram, chain$ とし、抽出したサンプルを学習して、以下の式（2）のパラメータ（回帰係数） $b_0, b_1, b_2, b_3$ を推定することを指す。

【0152】

【数2】

$$\text{Pr} = \frac{1}{1 + \exp(b_0 + b_1dist + b_2gram + b_3chain)} \quad \dots (2)$$

20

【0153】

抽出したサンプルから学習したパラメータ（回帰係数）は、例えば $b_0 = -1.425$ 、 $b_1 = -0.564$ 、 $b_2 = 11.036$ 、 $b_3 = 3.115$ と推定される（10000サンプルから回帰分析）。この場合、これらのパラメータを当てはめた式（3）が参照確率を求めるための回帰式である。

30

【0154】

【数3】

$$\text{Pr} = \frac{1}{1 + \exp(-1.425 - 0.564 \times dist + 11.036 \times gram + 3.115 \times chain)} \quad \dots (3)$$

40

【0155】

推定されるパラメータ（回帰係数） $b_0, b_1, b_2, b_3$ の値は、文書記憶手段2で記憶する文書データによって異なる。例えば、文書記憶手段2で記憶する文書データが書き言葉である新聞記事のみからなる場合と話し言葉である発話を文書データに変換したもののみからなる場合とでは、夫々推定されるパラメータは異なる。また、書き言葉として同種の新聞記事のみからなる文書データに対しても、その文書データの量、文書データの文書の内容によって推定されるパラメータの値 $b_0, b_1, b_2, b_3$ は異なる。そこで本発明では、話し言葉での回帰分析のために、書き言葉と話し言葉とで区別して文書データを記憶しておき、話し言葉からなる文書データに対しても回帰分析によってパラメータ

50

を推定し、参照確率を算出するための回帰式を記憶しておく。なお、受付装置4, 4, ...で受け付ける言葉が、音声入力された発話ではなく文字入力によって書き言葉からなる文章を入力したものに限定されている場合は、話し言葉と書き言葉とで文書データを区別せずに文書記憶手段2で記憶する構成としてもよい。

【0156】

以上の回帰分析により、式(3)の回帰式の特徴量  $d i s t, g r a m, c h a i n$  に対するパラメータが求められる。したがって、文単位検索装置1のCPU11が文単位の各単語の特徴量  $d i s t, g r a m, c h a i n$  からなる特徴パターンを特定することにより、当該特徴パターンを有する単語の参照確率を算出することができる。

【0157】

### 3-3. 文単位毎の顕現性の定量化

回帰分析により回帰式が得られたため、文単位検索装置1のCPU11は、文単位毎に抽出された単語毎に特徴量  $d i s t, g r a m, c h a i n$  を特定することにより、単語毎の参照確率を算出することができる。そこで、文単位検索装置1のCPU11は、文書記憶手段2で記憶しているタグ付け済みの文書データを取得して文毎に分別し、当該文で出現する単語又は参照する単語に対して特徴パターンを特定し参照確率を算出する。これにより、先行する文の文脈上の意味が反映された文毎の意味のまとまりを定量的に表すことができる。

【0158】

文単位検索装置1のCPU11が回帰分析後に、文書記憶手段2で記憶している文書データの文毎に、単語及び単語毎の参照確率(重み付き単語群)を算出する処理について以下に説明する。

【0159】

文単位検索装置1のCPU11は、文書記憶手段2で記憶している文書データを取得して、文書データに含まれる文毎にその文と先行の文とにおける各単語の文法的な特徴パターンを特定し、特定した特徴パターンと回帰式とに基づいて文毎に各単語の参照確率を算出して予め記憶する。

【0160】

文単位検索装置1のCPU11は、各単語と夫々の単語の参照確率との組(重み付き単語群)を各文単位毎に対応付けて記憶しておく。即ちCPU11は、文書集合から取得する全文書の全文について記憶する処理を行なう。一方、CPU11は、後の検索処理において、全文書の全文の内の、受け付けた言葉と文脈上の意味が類似する文を抽出する。したがって、この場合、全文書の全文を一つ一つ読み出して夫々に対応付けられている各文の文脈上の意味を表わす重み付き単語群を読み出すのでは処理の負荷が大きい。

【0161】

そこで、文単位検索装置1のCPU11は、各文に対して先行の文の文脈上の意味を表わした重み付き単語群を、後の処理で全文書の全文を一つ一つ読み出すことなしに抽出する処理を可能にするために、各文毎に算出した重み付き単語群をデータベース化して索引付けしておく処理を行なう。

【0162】

図9及び図10は、実施の形態1における文単位検索装置1のCPU11が、文書記憶手段2で記憶しているタグ付け済みの文書データの文毎に単語の参照確率を算出し、記憶する処理手順を示すフローチャートである。図9及び図10のフローチャートに示す処理は、文単位毎に、各単語に対して特定した特徴パターンと、特徴パターンに対応する回帰係数とを使用して参照確率を算出する処理、算出した参照確率を単語との組で予め記憶しておく処理に対応する。

【0163】

文単位検索装置1のCPU11は、文書記憶手段2から文書集合接続手段16を介してタグ付け済みの文書データを取得する(ステップS301)。CPU11は、取得した文書データに付加されたタグ < s u > を文字列解析によって識別して文に分別する(ステッ

10

20

30

40

50



プ S 3 0 2 )。次に CPU 1 1 は、文を示す < s u > 内の各タグを文字列解析によって識別し、文に対し、当該文で出現する単語又は参照される単語を抽出し (ステップ S 3 0 3 )、当該文書データについて参照確率の算出を行う間は、抽出した単語を一時記憶領域 1 4 で記憶する (ステップ S 3 0 4 )。

【 0 1 6 4 】

CPU 1 1 は、一時記憶領域 1 4 に記憶した、当該文を含む文書データについての単語に対し、単語に付加されたタグを文字列解析によって識別して *d i s t , g r a m , c h a i n* からなる特徴パターンを特定する (ステップ S 3 0 5 )。次に CPU 1 1 は、特定した特徴パターンの各特徴量を式 ( 3 ) に代入し参照確率を算出する (ステップ S 3 0 6 )。

10

【 0 1 6 5 】

CPU 1 1 は、文に対する各単語の参照確率を、一時記憶領域 1 4 で記憶している全単語に対して算出したか否かを判断する (ステップ S 3 0 7 )。CPU 1 1 が全単語に対して参照確率を算出していないと判断した場合は ( S 3 0 7 : N O )、CPU 1 1 は、処理をステップ S 3 0 5 に戻し、他の単語についての特徴パターンの特定及び参照確率の算出を継続する。一方、CPU 1 1 が全単語に対して参照確率を算出したと判断した場合は ( S 3 0 7 : Y E S )、CPU 1 1 は、一時記憶領域 1 4 で記憶している単語及び各単語に対して算出した参照確率の組 ( 重み付き単語群 ) を *s a l i e n c e* 属性を付加して記憶する (ステップ S 3 0 8 )。この際、CPU 1 1 は参照確率を所定の値で絞込み、参照確率が所定の値未満である単語については記憶しない。

20

【 0 1 6 6 】

次に、CPU 1 1 は、現在の文に対して付加した単語及び各単語の参照確率の組 ( 重み付き単語群 ) を後に抽出できるように、索引付けして重み付き単語群のデータベースに記憶する (ステップ S 3 0 9 )。CPU 1 1 はデータベースを記憶手段 1 3 に記憶してもよいし、文書集合接続手段 1 6 を介して文書記憶手段 2 に記憶してもよい。なお、CPU 1 1 は、索引付けの処理の 1 つとして以下のような処理を実行する。

【 0 1 6 7 】

CPU 1 1 は例えば、ステップ S 3 0 8 で得られた重み付き単語群の内の、一の単語の参照確率に注目し、一の単語の参照確率が所定値以上であるか否かを判定する。次に、CPU 1 1 は重み付き単語群の内の、他の一の単語の参照確率が所定値以上であるか否かを判定する。CPU 1 1 は、算出した重み付き単語群を、一の単語の参照確率が所定値以上のグループ、一の単語の参照確率が所定未満のグループのいずれに属するか、さらに一の単語の参照確率が所定値以上のグループに属する場合は、さらに他の単語の参照確率が所定値以上のグループ、他の単語の参照確率が所定未満のグループのいずれに属するかを判定しておく。CPU 1 1 は、このような処理を繰り返して算出した重み付き単語群がいずれのグループに属するかを判定し、属するグループの識別情報に対応付けて記憶しておく。この索引付けの処理は例えば、*k-d tree* 探索アルゴリズムを適用することができる。

30

【 0 1 6 8 】

CPU 1 1 は、ステップ S 3 0 1 で取得した文書データ中の全文について各文毎に重み付き単語群を対応付ける処理を終了したか否かを判断する (ステップ S 3 1 0 )。CPU 1 1 は、文書データ中の全文について各文毎に重み付き単語群を対応付ける処理を終了したか否かを以下のように判断する。例えば、現在の文を挟む < s u > < / s u > の後に、< s u > タグが後続するか否かを判断し、後続しないと判断した場合は終端であると判断することができる。CPU 1 1 がステップ S 3 0 1 で取得した文書データ中の全文について各文毎に重み付き単語群を対応付ける処理を終了していないと判断した場合は ( S 3 1 0 : N O )、CPU 1 1 は、処理をステップ S 3 0 2 に戻し、次の文に対して処理を継続する。一方、CPU 1 1 がステップ S 3 0 1 で取得した文書データ中の全文について各文毎に重み付き単語群を対応付ける処理を終了したと判断した場合は ( S 3 1 0 : Y E S )、CPU 1 1 は、文書データで抽出されて一時記憶領域 1 4 に記憶していた単語を消去する (ステップ S 3 1 1 )。

40

50

## 【0169】

CPU11は、全文書データについて、単語及び単語の参照確率を *salience* 属性によって記憶する処理を終了したか否かを判断する(ステップS312)。CPU11が全文書データについて、単語及び単語の参照確率を *salience* 属性によって記憶する処理を終了していないと判断した場合は(S312:NO)、CPU11は、処理をステップS301へ戻し、別の文書データを取得して処理を継続する。CPU11が全文書データについて、単語及び単語の参照確率を *salience* 属性によって記憶する処理を終了したと判断した場合は(S312:YES)、CPU11は、単語の参照確率を算出して予め記憶する処理を終了する。

## 【0170】

次に、文単位検索装置1のCPU11が図9及び図10のフローチャートに示した処理を図5に示した文書データに対して行った場合について具体的に説明する。

## 【0171】

図11は、実施の形態1における文単位検索装置1のCPU11が、文書データに示される文書を文毎に分別した一例を示す説明図である。

## 【0172】

文単位検索装置1のCPU11は、ステップS301及びステップS302の処理により、文書記憶手段2で記憶している文書データから、<su>タグを識別して文毎に分別する。図11に示す例では、文は $s_1$ 「祭とは、神霊などを祀る儀式。」、 $s_2$ 「祭礼、祭祀とも呼ばれる。」、 $s_3$ 「九州地方北部では、秋に行われるものに対して(お)くんちと称する場合もある。」に分別される。文単位検索装置1のCPU11によるステップS303の処理により、文 $s_1$ 、 $s_2$ 、 $s_3$ から抽出される単語は、単語のリストに記憶された単語と一致する「祭」、「神霊」、「儀式」、「祭礼」、「祭祀」、「九州」、「九州地方」、「九州地方北部」、「秋」、「くんち」、「場合」である(図6参照)。

## 【0173】

文単位検索装置1のCPU11は、ステップS305の処理により、各単語群の文 $s_3$ での顕現性(参照確率)を定量的に求めるために、各単語群の特徴量 *dist*, *gram*, *chain* からなる特徴パターンを特定する。例えば、文 $s_3$ での「九州」(識別番号:9714)(図6参照)の特徴パターンは以下のように特定される。

## 【0174】

図11の説明図に示すように、文 $s_3$ での「九州」の *dist* は、最近出現した文 $s_3$ と、後続の文 $s_4$ との距離1により *dist* = 1である。また、文 $s_3$ での「九州」の *gram* は、最近「九州」が出現した文 $s_3$ では「九州」に係るのは助詞ではなく「地方」へ係るために名詞接続と特定でき *gram* = 名詞接続である。文 $s_3$ での「九州」の *chain* は、 $s_1$  から  $s_3$  まで「九州」が出現した回数は一回であるので *chain* = 1である。したがって、特徴パターン  $f(s_3, \text{九州}) = (\text{dist} = 1, \text{gram} = \text{名詞接続}, \text{chain} = 1)$  と特定される。したがって、文単位検索装置1のCPU11は、図9及び図10のフローチャートのステップS306の処理により、式(3)に特徴量 *dist*, *gram*, *chain* の値を代入して参照確率を算出する。

## 【0175】

ここで、*gram* で表される特徴量の代入値は、文書記憶手段2で記憶した文書データからサンプル( $s, w$ )を抽出し、夫々に対して算出した単語  $w$  の参照確率を *gram* 毎に平均値を算出し代入値とする。例えば、抽出したサンプル( $s, w$ )のうち、*gram* = 八を有する単語に対して算出した参照確率の平均値が特徴量 *gram* が「八」である場合に代入する値である。実施の形態1では、例として、*gram* = 八の場合は *gram* = 0.0540、*gram* = ガの場合は *gram* = 0.0288、*gram* = ノの場合は *gram* = 0.0198、*gram* = ヲの場合は *gram* = 0.0179、*gram* = ニである場合は *gram* = 0.0124、*gram* = 名詞接続である場合は、*gram* = 0.00352が算出される。

## 【0176】

10

20

30

40

50

なお、単語が、助詞「ハ」に係る場合、助詞「ガ」に係る場合、助詞「ノ」に係る場合、助詞「ヲ」に係る場合での、当該単語が後続の文で出現する参照確率の平均値は、「ハ」(主題)「ガ」(主語)「ノ」「ヲ」(目的語)の順に高く、当該文での中心であるか否かを示す中心化理論で定式化している主題>主語>目的語...の序列とほぼ整合する。

【0177】

文s<sub>3</sub>での「九州」の参照確率(文s<sub>4</sub>で「九州」が出現又は参照される確率)は、特定した特徴量に基づいて以下式(4)のように算出される。

【0178】

【数4】

$$Pr = \frac{1}{1 + \exp(-1.425 - 0.564 \times 1 + 11.036 \times 0.00352 + 3.115 \times 1)}$$

$$= 0.238$$

... (4)

10

20

【0179】

式(4)に示したように、文s<sub>3</sub>での「九州」の参照確率は0.238と算出される。算出された参照確率は文s<sub>3</sub>に対して記憶される。文単位検索装置1のCPU11は、文s<sub>3</sub>に対し単語をリストで記憶した識別番号で表し、参照確率を対応付けて記憶する。本発明では、文の単位を区切る<su>タグに対して属性名salienceを定義し、属性値は単語の識別番号及び参照確率の組を羅列したものと定義して以下のように文毎に単語及び該単語の参照確率(重み付き単語群)を記憶する。

【0180】

<su salience = " 単語<sub>1</sub> の識別番号 : 単語<sub>1</sub> の参照確率 単語<sub>2</sub> の識別番号 : 単語<sub>2</sub> の参照確率 単語<sub>3</sub> の識別番号 : 単語<sub>3</sub> の参照確率 ... " > ... </su>

30

【0181】

図12は、実施の形態1における文単位検索装置1のCPU11が、参照確率を算出した結果を付与して文書記憶手段2に記憶させる文書データの一例を示す説明図である。文s<sub>3</sub>では「九州」(9714)の参照確率(文s<sub>3</sub>での重み値。以下同様)が0.238、「九州地方北部」(9716)の参照確率が0.1159、...と記憶され、後続の文s<sub>4</sub>では「九州」(9714)の参照確率が0.238、「祭」(22953)の参照確率が0.1836、...と記憶される。文毎に異なる単語及び参照確率の組(重み付き単語群)が記憶され、文毎の意味のまとまりを表す情報として検索に使用することができる。文s<sub>3</sub>及び文s<sub>4</sub>で、「九州」(9716)は、同値の参照確率が算出されているが、文s<sub>5</sub>, 文s<sub>6</sub>, ...と続く毎に、九州地方に限らない「祭」についての記述が続く場合は「九州」の参照確率は次第に低下していくと考えられる。

40

【0182】

図13は、実施の形態1における文単位検索装置1のCPU11が、文単位毎に算出した重み付き単語群を索引付けして記憶した場合のデータベースの内容例を示す説明図である。なお、図13の内容例は、図12の内容例に示した文s<sub>4</sub>に対応付けられる重み付き単語群が、図9及び図10のフローチャートに示したCPU11のステップS309によって索引付けされたデータに相当する。

【0183】

図13に示すように、CPU11は重み付き単語群を、いずれのグループに属するかを示す情報(k-d tree ノードID)に対応付けて記憶しておく。さらにその際、CPU11

50

は、その重み付き単語群がいずれの文書データの文単位に対応付けられているかを特定できるように、タグ付け済み文書データのファイル名及び文書データ中の位置（タグ情報）を記憶しておく。これにより、後の処理で受け付けた言葉に対して求めた重み付き単語群と類似する重み付き単語群が対応付けられている文単位を抽出することが容易になる。

【0184】

図14は、文単位検索装置1のCPU11により文毎に記憶される単語及び該単語に対して算出された参照確率の組が、文が続くにつれてどのように変化するかを示す説明図である。図14では、文 $s_1$ 、文 $s_2$ 、文 $s_3$ 、文 $s_4$ と続くにつれて、時系列で文脈が動的に変化することに応じて、夫々の文で顕現性の高い単語が夫々異なることが判る。

【0185】

#### 4. 検索処理

##### 4-1. ユーザから入力された言葉の受け付け

次に、実施の形態1における検索処理について説明する。検索処理は、受付装置4, 4, ...でユーザから入力されるキーワード又は音声等の言葉を受け付けたことを起点として開始する。

【0186】

受付装置4のCPU41は、操作手段45を介してユーザが入力する文字列を検知して一時記憶領域44に記憶する処理、又は音声入出力手段47を介してユーザが入力する音声を検知して文字列に変換し一時記憶領域44に記憶する処理が可能である。また、受付装置4のCPU41はユーザが入力する文字列を解析して一文一文に分別する機能を有する。例えば、日本語の場合は句点「。」、英語の場合はピリオド「.」等の所定の文字を識別して分別するのでもよい。また、Enterキーが押下されたことを操作手段45を介して検知する都度、Enterキーが入力されるまでの文字列を一文と分別するのでもよい。ユーザからの音声入力に対しては、例えば、音声認識機能によって音声を文字列に変換し、変換した文字列から文字列解析によって文に分別してもよいし、無音を検出したところで文に分別してもよい。受付装置4のCPU41は、分別した一文一文をテキストデータとして通信手段48を介して文単位検索装置1へ送信する。

【0187】

##### 4-2. 受け付けた言葉に対する意味のまとまりの定量化

次に、文単位検索装置1のCPU11が、受付装置4, 4, ...で受け付けた言葉を示すテキストデータを受信した場合に、文書記憶手段2で記憶している文書中の文を検索する処理について説明する。受け付けた言葉を示すテキストデータに対しても、意味のまとまりの定量化、即ち当該テキストデータの単語抽出及び単語の参照確率の算出を行う。これにより、ユーザが言葉を入力するときにユーザの潜在的な意識にある先行の言葉からの流れに応じた文脈を反映した意味のまとまりを表わす情報を、後述する検索処理における検索要求として自動的に作成することができる。

【0188】

文単位検索装置1のCPU11は、ユーザから受け付けた言葉を示すテキストデータをパケット交換網3及び通信手段15を介して受付装置4, 4, ...から受信した場合、一時記憶領域14に受信した順にテキストデータを記憶すると共に、受信したテキストデータで示される文に対して形態素解析及び統語解析を行う。また、受信したテキストデータで示された文 $s$ と、文 $s$ より以前に受信したテキストデータで示された文に出現した単語 $w$ との対 $(s, w)$ に対し、特徴量 $dist, gram, chain$ で表される特徴パターン $f(s, w)$ を特定する。

【0189】

文単位検索装置1のCPU11は、受信したテキストデータの文 $s$ での単語 $w$ の特徴パターン $f(s, w)$ を特定した場合、特定した特徴パターンと先に得られた回帰式とに基づいて参照確率を算出する。文単位検索装置1のCPU11は、各単語について参照確率を算出し、各単語と各単語について算出した参照確率とを用いて、既に文単位に対応付けて記憶してある重み付き単語群、即ち各単語と各単語の参照確率との組と比較する処理を

10

20

30

40

50

おこなって文単位の検索を行う。

【0190】

なお、文単位検索装置1のCPU11は、受付装置4, 4, ...からテキストデータのみならず、ユーザから入力された発話の音声データも受信することが可能である。この場合、音声データをテキストデータと同様に音声データに表わされている単語の文法上の特徴パターンを特定することにより、同様の処理を行なう。また、音声データの場合は音声データで得られる特徴を、その単語の顕現性が高いか否かを判断するための特徴量として扱うことも可能である。例えば、CPU11は、単語が出現又は参照された場合に、先行の言葉で出現又は参照されてからの時間差を一つの特徴量として扱うことができる。またCPU11は、その単語が出現又は参照された直近の先行の言葉中で、その単語が発声されたときの発話速度及び/又は音声の周波数を他の特徴量として扱うことができる。これらは、テキストデータに変換された後では検知することができない、時間情報又は単語にこめられた感情を定量的に表わす情報である。

10

【0191】

受付装置4がユーザから入力された言葉を受け付けて文単位検索装置1へ送信し、文単位検索装置1のCPU11が受付装置4から受信したテキストデータに基づいて文書記憶手段2で記憶している文書データから検索を行う処理手順についてフローチャートを用いて説明する。図15、図16、及び図17は、実施の形態1における文単位検索装置1及び受付装置4の検索処理の処理手順を示すフローチャートである。

【0192】

受付装置4のCPU41は、ユーザによる文字列入力操作を操作手段45を介して検知したか否か、又はユーザによる音声入力を音声入出力手段47を介して検知したか否かを判断する(ステップS401)。CPU41がユーザによる文字列入力操作又は音声入力を検知していないと判断した場合は(S401:NO)、CPU41は、処理をステップS401へ戻し、ユーザによる文字列入力操作又は音声入力を検知するまで待機する。

20

【0193】

一方、受付装置4のCPU41がユーザによる文字列入力操作又は音声入力を検知したと判断した場合は(S401:YES)、受付装置4のCPU41は、入力された文字列又は音声入力を変換した文字列から、入力された言葉を一文に分別して一時記憶領域44に記憶し(ステップS402)、ユーザから入力された言葉をパケット交換網3を介して文単位検索装置1へ送信する(ステップS403)。

30

【0194】

文単位検索装置1のCPU11は、受付装置4から、ユーザによって入力された言葉を受信し(ステップS404)、CPU11は、受信した言葉を文として一時記憶領域14に受信順にテキストデータで記憶する(ステップS405)。このとき、テキストデータ毎に文識別番号を付加して記憶してもよい。

【0195】

CPU11は、記憶したテキストデータを形態素解析及び統語解析し(ステップS406)、解析によって抽出された単語を一時記憶領域14に記憶する(ステップS407)。このときCPU11は、リストに記憶してある単語と照合し、リストの識別番号で単語を記憶する。

40

【0196】

なお、文単位検索装置1のステップS407における処理により、一時記憶領域14には、一連として入力された言葉(発話)の中で一度は出現又は参照された単語が記憶されることになる。なお、ステップS407における単語の抽出は必ずしも行わなくてもよい。その場合は、リストに記憶してある全単語に対し、後述する特徴パターンの特定の処理を行う。

【0197】

CPU11は、一時記憶領域14に記憶している単語夫々に対し、過去に受信して記憶してあるテキストデータ及びステップS406の形態素解析及び統語解析の結果に基づい

50

て、特徴パターンを特定する(ステップS408)。CPU11は、特定した特徴パターンの特徴量を、予め話し言葉について回帰分析して求めた参照確率を算出するための回帰式に代入し、単語毎に参照確率を算出する(ステップS409)。CPU11は、一時記憶領域14で記憶している全単語について参照確率を算出したか否かを判断する(ステップS410)。CPU11が記憶している全単語について参照確率を算出していないと判断した場合は(S410:NO)、処理をステップS408へ戻し、別の単語について特徴パターンの特定及び参照確率の算出の処理を行う。

**【0198】**

CPU11が記憶している全単語について参照確率を算出したと判断した場合は(S410:YES)、一時記憶領域14に夫々参照確率を算出して記憶している全単語に対し、所定値以上の参照確率が算出された単語に絞り込む(ステップS411)。参照確率が極端に低い単語を除去することにより、後の演算によるCPU11自身への負荷を低減させるためである。CPU11は、受け付けた言葉に対して絞り込まれた単語及び単語の参照確率に基づいて以下のような検索処理を行う。

10

**【0199】**

これまでの処理により、受け付けた言葉に対し、以前に受け付けた言葉から続く流れ上の意味のまとまりを定量的に表わす単語と単語の参照確率の組(重み付き単語群)を検索要求として生成することができた。以下の検索処理(一点鎖線で囲まれたステップS412からステップS416まで)は、受け付けた言葉に対して得られた重み付き単語群と、予め記憶してある文単位の重み付き単語群とを比較し、夫々の重み付き単語群の内の複数の単語の重み値の分布が類似するか否かによって、言葉と文とで意味が類似するか否かを判定し、類似する文を抽出する処理の一例である。

20

**【0200】**

CPU11は、記憶手段13又は文書記憶手段2のデータベースから、各文に対応付けられて記憶されている単語と単語の参照確率との組(以下重み付き単語群という)を読み出す(ステップS412)。

**【0201】**

このとき、CPU11は、ある程度類似する重み付き単語群を絞り込んで読み出すことができるように、ステップS411までの処理で得られた受け付けた言葉に対応付けられる重み付き単語群が、データベースに記憶してある重み付き単語群同様にいずれのグループに属するかを判定する。CPU11は、受け付けた言葉に対応付けられた重み付き単語群が属するグループの重み付き単語群をデータベースから読み出す。これにより、全く類似しない重み付き単語群と比較することを回避し、ある程度類似する重み付き単語群を絞り込んで抽出することができる。

30

**【0202】**

次にCPU11は、ステップS412で読み出した重み付き単語群から、受け付けた言葉の重み付き単語群と同一の単語を含む重み付き単語群を抽出する(ステップS413)。CPU11は、抽出した文と同一の単語夫々について、参照確率の差分を算出する(ステップS414)。CPU11は、同一の単語の数の多い順及び同一の単語の参照確率の差分が小さい順に、抽出した重み付き単語群に類似度を付与し(ステップS415)、抽出した重み付き単語群が対応付けられている文を文書集合の文書データから読み出す(ステップS416)。このとき、CPU11は、類似度が所定値以上の重み付き単語群のみに対応する文を読み出してもよい。CPU11は、抽出した文を類似度でソートする(ステップS417)。

40

**【0203】**

上述のステップS412からステップS417までの処理により、受け付けた言葉に対して得られた重み付き単語群の内の複数の単語の重み値の分布と、類似する重み値の分布を有する重み付き単語群が対応付けられた文を抽出することができる。

**【0204】**

次にCPU11は、各文を表すテキストデータを検索結果のテキストデータとして受付

50

装置 4 へ通信手段 15 を介して送信する (ステップ S 4 1 8)。

【 0 2 0 5 】

受付装置 4 の CPU 4 1 は、検索結果のテキストデータを通信手段 4 8 を介して受信し (ステップ S 4 1 9)、受信したテキストデータを表示手段 4 6 を介してモニタ等に表示し (ステップ S 4 2 0)、処理を終了する。

【 0 2 0 6 】

受付装置 4 の CPU 4 1 は、ユーザからの言葉の入力を検知する都度、一文に分別したテキストデータ又は音声データを文単位検索装置 1 へ送信する。文単位検索装置 1 の CPU 1 1 は、受付装置 4 からテキストデータ又は音声データ、音声データと共に送信される情報を受信する都度、単語及び単語毎の参照確率を算出して、ユーザから受け付けた言葉 10 に対し、先行の言葉からの流れが反映された意味のまとまりを表わす情報、即ち重み付き単語群を検索要求として作成する。文単位検索装置 1 の CPU 1 1 は、受け付けた言葉に対して作成した検索要求 (重み付き単語群) に基づいて記憶している文書データから文単位を抽出し、検索結果としてテキストデータを送信する。

【 0 2 0 7 】

実施の形態 1 における受付装置 4 の CPU 4 1 は、検索結果のテキストデータを受信する都度、モニタ等に表示する。したがって、受付装置 4 ではユーザから言葉が入力される都度、当該言葉と意味のまとまりが類似するテキストデータが検索結果として表示される。

【 0 2 0 8 】

なお、受付装置 4 は、必ずしもユーザから言葉が入力される都度毎回テキストデータを送信し、検索結果を受け付けて表示する構成としなくともよい。例えば、所定の期間中に入力された複数の言葉に相当するテキストデータ又は音声データを文単位検索装置 1 へ送信し、複数の言葉に対応する検索結果を受け付けて表示する構成でもよい。

【 0 2 0 9 】

図 1 5、図 1 6 及び図 1 7 のフローチャートに示した文単位検索装置 1 の CPU 1 1 による処理の詳細を具体例を挙げて以下に説明する。

【 0 2 1 0 】

図 1 8 は、実施の形態 1 における文単位検索装置 1 の CPU 1 1 が、受付装置 4 から受信したテキストデータに対して特定した特徴パターンの例を示す説明図である。図 1 8 中の文単位  $S_{i-2}$ 、文単位  $S_{i-1}$ 、文単位  $S_i$  は夫々、受信した各テキストデータで示される文である。

【 0 2 1 1 】

図 1 8 中の文単位  $S_i$  での、当該文単位  $s_i$  及び先行する文単位に含まれる単語「おくんち」とのサンプル ( $s_i$ 、おくんち) の特徴パターンは以下のようにして特定される。現在の文  $s_i$  及び先行する文のうち、単語「おくんち」が最近出現又は参照された文  $s_{i-2}$  との距離の特徴量 ( $dist$ ) は、 $dist = 3$  である。また、単語「おくんち」が最近出現又は参照された  $s_{i-2}$  での「おくんち」が係っている格助詞は「って」であるため、 $gram = \text{ッテ}$  である。更に、文  $s_i$  より先行の文  $s_{i-2}$  で単語「おくんち」が出現又は参照されたため  $chain = 1$  である。したがって、特徴パターンは  $f(s_i, \text{おくんち}) = (dist = 3, gram = \text{ッテ}, chain = 1)$  と特定される。英語の場合、 $gram$  は前置詞によって特定される。

【 0 2 1 2 】

文単位検索装置 1 では、話し言葉についても文書記憶手段 2 で記憶している文書データについて回帰分析を行い、特徴パターンを特定した場合に特徴量を代入することで参照確率を算出することができる回帰式が予め導出されている。したがって、文単位検索装置 1 の CPU 1 1 は、文  $s_i$  の「おくんち」に対して、特定した特徴パターンの特徴量  $dist$ ,  $gram$ ,  $chain$  に基づいて参照確率を算出することができる。更に、文単位検索装置 1 の CPU 1 1 は、文  $s_i$  について過去に出現又は参照された単語も含めて参照確率を算出し、単語と単語の参照確率とを求める。文単位検索装置 1 の CPU 1 1 は、求め 50

た単語と参照確率とに基づいて、文書記憶手段2で記憶してある *s a l i e n c e* 属性を予め記憶してある文単位から同一の単語の参照確率が所定の値以上である文単位を直接的に抽出する。文単位検索装置1のCPU11は、抽出した文を示すテキストデータを通信手段15を介して受付装置4へ送信する。

#### 【0213】

このような文単位検索装置1のCPU11の処理により、受信したテキストデータが表示する言葉の意味のまとまりを当該言葉毎に単語及び単語の参照確率(重み値)で表すことができる。また、予め文書記憶手段2で記憶してある文書データの各文についても、意味のまとまりを表す単語及び単語の参照確率(重み付き単語群)が記憶されるので、ユーザから受け付けた言葉に対し、抽出された単語の参照確率が類似するか否かによって意味のまとまりが類似する文を直接的に検索することができる。

10

#### 【0214】

(実施の形態2)

実施の形態2では、事前処理の段階で文書記憶手段2で記憶した文書データの文毎に、抽出した単語と単語毎に算出した参照確率との組(重み付き単語群)を顕現性ベクトルとして扱う。さらに、受け付けた言葉に対して算出する単語と単語毎に算出した参照確率との組(重み付き単語群)も顕現性ベクトルとして扱う。そして検索処理の段階においては、実施の形態1に示したように、受け付けた言葉の重み付き単語群の内の複数の単語の重み値の分布と、予め文毎に対応付けてある重み付き単語群の内の複数の単語の重み値の分布とが類似する条件にあるか否かを、同一の単語が記憶されており、同一の単語の差分が小さいか否かで判断した。これに対し、実施の形態2では、夫々の重み付き単語群を顕現性ベクトルで表わし、類似する条件にあるか否かを顕現性ベクトル間の距離の長さによって判断する。

20

#### 【0215】

実施の形態2における、本発明に係る文単位検索装置1を用いた検索システムの「1.ハードウェアの構成及び概要」、及び「2.文書データの取得及び自然言語解析」については、実施の形態1と同様であるため説明を省略する。「3.文書データの文毎の意味のまとまりの定量化」、及び「4.検索処理」について以下に説明するが、実施の形態1と同一の符号を用いて説明する。なお、「3.文書データの文毎の意味のまとまりの定量化」、及び「4.検索処理」についても、実施の形態1と共通する点については詳細な説明を省略する。

30

#### 【0216】

3.文書データの文毎の意味のまとまりの定量化

3-1.文毎の意味のまとまりの定義

実施の形態2では、実施の形態1と同様に文毎の意味のまとまりを定量的に表す情報は、ユーザが当該文を使用(発話、筆記、聴取又は読解)するときに、ユーザが注目している単語群と、ユーザが各単語に注目する度合い、即ち顕現性(*s a l i e n c e*)を定量的に示す値(単語の重み値)とで表す。また、実施の形態1と同様に、顕現性を定量的に示す重み値として後続の文で出現する又は参照される確率を示す参照確率を使用する。

#### 【0217】

3-2.回帰モデル学習

実施の形態2でも、参照確率については実施の形態1の3-2.回帰モデル学習と同様に、文書記憶手段2で記憶している文書データのサンプルに対する回帰分析によって得られる回帰係数を含む回帰式を用いて算出する。

40

#### 【0218】

3-3.文単位毎の顕現性の定量化

実施の形態2でも、文単位検索装置1のCPU11は、回帰分析によって得られた回帰係数を含む回帰式を使用して、抽出された単語毎に特徴量 *d i s t , g r a m , c h a i n* を特定することで単語毎の参照確率を算出することができる。ここで、単語毎の参照確率をその単語の重み値として付与した重み付き単語群が得られる。実施の形態2では、文

50



毎の意味のまとまりを表わす重み付き単語群は、単語を夫々次元とし、単語毎に算出した参照確率を各単語に対応する次元成分の要素として持つ顕現性ベクトルとして扱う。つまり、文書記憶手段2で記憶される文書データ中の文の意味のまとまりは、文書記憶手段2で記憶される文書データから抽出し、図6に示すリストに記憶している31245次元の多次元空間におけるベクトルで表すことができる。

【0219】

したがって、(あい, あいだ, あいまい, ..., Z, Zくん)という単語群からなる31245次元の基底空間に対し、図11に示した文 $s_3$ の顕現性ベクトル $v(s_3)$ は、文 $s_3$ での9714番目の「九州」次元に対応する要素が参照確率の大きさ(重み値)0.238で表され、また、9716番目の「九州地方北部」次元に対応する要素が参照確率の大きさ0.1159で表されるので、(0, 0, ..., 0.238, 0, 0.1159, ..., 0)と31245次元のベクトルで表現して扱うことができる。

10

【0220】

なお、実施の形態2において文単位検索装置1のCPU11が参照確率を算出した結果を付与して文書記憶手段2に記憶させる文書データは、実施の形態1の図11の説明図に示した文書データと同様である。即ち、文書記憶手段2に記憶される文書データには、次元の番号及び次元成分の要素である参照確率の値が記憶される。実施の形態2における文単位検索装置1のCPU11が、文書記憶手段2で記憶しているタグ付け済みの文書データの文毎に単語の参照確率を算出し、文毎に対応付けてデータベースに記憶する処理手順は、実施の形態1と同様であるため説明を省く。

20

【0221】

#### 4. 検索処理

次に、実施の形態2における検索処理について説明する。「4-1. ユーザから入力された言葉の受け付け」については、受付装置4のCPU41が行う処理については実施の形態1と同様である。

【0222】

#### 4-2. 受け付けた言葉に対する意味のまとまりの定量化

文単位検索装置1のCPU11が、受付装置4で受け付けた言葉を示すテキストデータを受信した場合に、文書記憶手段2で記憶している文書中の文を検索する処理について説明する。文単位検索装置1のCPU11は、受け付けた言葉を示すテキストデータに対しても、受け付けた言葉の文脈上の意味のまとまりを単語の多次元空間における方向性を示す顕現性ベクトルで表す。

30

【0223】

文単位検索装置1のCPU11は、実施の形態1での処理同様に、受付装置4から受信したテキストデータに対してリストに記憶された31245次元の単語に対する特徴量  $dist, gram, chain$  で表される特徴パターンを特定する。なお、過去に一連として受信したテキストデータで出現していない単語については、対応する次元成分の要素を0として特徴パターンの特定を省く。

【0224】

特徴パターンを表す特徴量  $dist, gram, chain$  から、回帰式に基づいて次元成分の要素としての参照確率を夫々算出することができる。したがって、文単位検索装置1のCPU11は、テキストデータを受信する都度、受信したテキストデータで示される言葉のそれまでの文脈上の意味のまとまりを表わす顕現性ベクトルを算出することができる。

40

【0225】

文単位検索装置1のCPU11は、受け付けた言葉に対して算出した顕現性ベクトルと、文書記憶手段2で記憶してある、 $salience$ 属性を予め付加した文の顕現性ベクトルとの距離をベクトル演算によって直接算出し、距離が短い文を抽出する。図6の各単語を1次元とした場合の31245次元の多次元空間の中で意味のまとまりの方向性が類似する文を検索することができる。文単位検索装置1のCPU11は、抽出した文を示す

50

テキストデータを、通信手段 15 を介して受付装置 4 へ送信する。ベクトル演算を扱うことが可能なコンピュータを用いる場合は、文毎の意味のまとまりを顕現性ベクトルで表して直接的に演算をすることができる。

【0226】

文単位検索装置 1 の CPU 11 が、受付装置 4 で検索要求の言葉を示すテキストデータを受信し、受信したテキストデータに基づいて文書記憶手段 2 で記憶している文書データから顕現性ベクトルを用いて検索を行う処理手順について説明する。図 19 は、実施の形態 2 における文単位検索装置 1 及び受付装置 4 の検索処理の処理手順を示すフローチャートである。なお、図 19 のフローチャートに示す処理手順では、実施の形態 1 における図 15、図 16 及び図 17 のフローチャートに示した検索処理の処理手順と同一の処理につ

10

【0227】

図 19 のフローチャートに示す処理手順の内、一点鎖線で囲まれた各ステップ S501 からステップ S506 までの処理が、実施の形態 1 における図 15、図 16 及び図 17 のフローチャートに示した処理手順と異なる。実施の形態 1 におけるステップ S412 からステップ S416 までの処理の代わりに、実施の形態 2 における文単位検索装置 1 の CPU 11 により実行されるステップ S501 からステップ S506 までの処理について、以下に説明する。

【0228】

文単位検索装置 1 の CPU 11 は、一時記憶領域 14 に夫々参照確率を算出して記憶している全単語に対し、所定値以上の参照確率が算出された単語に絞り込み（ステップ S411）、絞り込まれた各単語と、算出された各単語の参照確率とに基づいて受け付けた言葉の顕現性ベクトルを算出する（ステップ S501）。

20

【0229】

ステップ S501 までの処理により、受け付けた言葉に対し、以前に受け付けた言葉から続く流れ上の意味のまとまりを定量的に表わす顕現性ベクトルを検索要求として生成することができた。以下の処理は、受け付けた言葉に対して得られた顕現性ベクトルと、予め記憶してある文毎の顕現性ベクトルとを比較し、夫々の顕現性ベクトルが表わす各単語の重み値の分布が類似するか否かを判定する処理の一例である。

【0230】

CPU 11 は、データベースに記憶してある重み付き単語群即ち顕現性ベクトルを読み出す（ステップ S502）。このとき、ステップ S411 までの処理で得られた受け付けた言葉に対応付けられる顕現性ベクトルが、データベースに記憶してある顕現性ベクトル同様にいずれのグループに属するかを判定する。CPU 11 は、受け付けた言葉に対応付けられた顕現性ベクトルが属するグループの顕現性ベクトルをデータベースから読み出す。これにより、各単語の重み値の分布が類似する顕現性ベクトルをある程度絞り込んで抽出することができる。

30

【0231】

CPU 11 は、受け付けた言葉に対応付けた顕現性ベクトルと読み出した顕現性ベクトルとの距離を算出する（ステップ S503）。CPU 11 は、読み出した顕現性ベクトルを、算出した距離が所定値未満である顕現性ベクトルに絞り込み（ステップ S504）、絞り込まれた顕現性ベクトルが対応付けられて記憶されている文を読み出す（ステップ S505）。CPU 11 は、読み出した文に算出した距離が短い順に類似度を付与する（ステップ S506）。

40

【0232】

実施の形態 2 における文単位検索装置 1 の CPU 11 によるステップ S501 からステップ S506 までの処理により、受け付けた言葉と文脈上の意味合いが類似する文が抽出される。

【0233】

その後の抽出された文に対するステップ S417 以降の処理は実施の形態 1 と同様であ

50

る。

【 0 2 3 4 】

なお、上述の処理手順の内の、CPU 11 が受け付けた言葉に対応付けた顕現性ベクトルと、読み出した顕現性ベクトルとの距離を算出するステップ S 5 0 3 の処理は、具体的には以下のように算出する。受け付けた言葉  $u_i$  に対応付けた顕現性ベクトルが  $v(u_i)$  と表わされ、読み出した顕現性ベクトルが  $v(s_i)$  と表わされる場合、CPU 11 は以下に示す式 (5) のように、コサイン距離を算出する。

【 0 2 3 5 】

【数 5】

$$\frac{\vec{v}(s_i) \cdot \vec{v}(u_i)}{|\vec{v}(s_i)| |\vec{v}(u_i)|} = \frac{\sum_{k=1}^{31245} \Pr(w_k | \text{pre}(s_i)) \Pr(w_k | \text{pre}(u_i))}{\sqrt{\sum_{k=1}^{31245} \Pr(w_k | \text{pre}(s_i))^2} \sqrt{\sum_{k=1}^{31245} \Pr(w_k | \text{pre}(u_i))^2}} \dots (5)$$

10

20

【 0 2 3 6 】

ただし、式 (5) に示したように距離を算出した場合、言葉の顕現性ベクトル  $v(u_i)$  と、読み出した顕現性ベクトル  $v(s_i)$  とが近いほど、算出したコサイン距離の値は大きくなる。したがって、CPU 11 はステップ S 5 0 6 において、算出したコサイン距離が大きい順に類似度を付与する。

【 0 2 3 7 】

このような文単位検索装置 1 の CPU 11 及び受付装置 4 の CPU 4 1 の処理により、受け付けた言葉の意味のまとまりを、当該言葉毎に各単語の参照確率を要素とした顕現性ベクトルで表すことができる。また、予め文書記憶手段 2 で記憶してある文書データの各文についても、意味のまとまりを表す各単語の参照確率を要素とした顕現性ベクトルが記憶してあるため、単語の多次元空間での方向性を表す顕現性ベクトル間の距離によって、意味のまとまりが類似する文を直接的に検索することができる。

30

【 0 2 3 8 】

(実施の形態 3)

実施の形態 1 又は 2 では、事前処理の段階の「3. 文書データの文単位毎の意味のまとまりの定量化」を行なう処理の中で、重み付き単語群として当該単語と単語の参照確率との組、又は顕現性ベクトルを文単位毎に対応付けて記憶しておいた。また、その後の「4. 検索の処理」でも「4-2. 受け付けた言葉に対する意味のまとまりの定量化」の処理の中で、重み付き単語群として単語と単語の参照確率との組、又は顕現性ベクトルを求めて受け付けた言葉に対応付けた。これに対し、実施の形態 3 では、文単位又は言葉毎に対応付けた重み付き単語群 (単語と単語の参照確率との組、又は顕現性ベクトル) に対し、各単語の顕現性を表わす重み値を、単語に関連の深い他の単語からの連想を加味して算出し直す処理を実行する。

40

【 0 2 3 9 】

具体的に連想とは、文単位毎に対応付けられている重み付き単語群の内のある単語が、その文単位又は先行の文単位に出現していない場合であっても、その単語と関連の深い単語の顕現性が高い場合はその単語もその文単位で注目されているはずであることをいう。

50

したがって、一の単語が注目されている時に同時に注目されやすい単語を関連語とする。そして、各単語の顕現性を表わす重み値に、関連の深い単語の顕現性からの影響を反映させる。

#### 【0240】

図20は、実施の形態3における本発明の検索方法に関わる、一の単語と関連の深い単語の顕現性の影響の概要を示す説明図である。図20の説明図は、一又は複数のユーザ間の会話の例を表わしている。会話は発話 $U_1$ 、 $U_2$ 、 $U_3$ 、 $U_4$ の集合であり、 $U_1$ 、 $U_2$ 、 $U_3$ 、 $U_4$ の順になされている。

#### 【0241】

ここで、発話 $U_1$ 、 $U_2$ 、 $U_3$ 、 $U_4$ にはいずれにも「大阪」は出現していない。また、 $U_1$ よりも先行の発話で「大阪」が出現しており、発話 $U_1$ 、 $U_2$ 、 $U_3$ 、 $U_4$ 夫々での「大阪」の顕現性がゼロではなく、ある程度の高さを有していたとしても、その後「大阪」は出現していないので、発話 $U_4$ の時点で「大阪」の顕現性を現す参照確率を定量的に算出した場合、その値が低下している可能性がある。

#### 【0242】

しかしながら、「大阪」という単語がそれまでの文単位又は言葉に出現していない場合であっても、発話 $U_1$ 、 $U_3$ には単語「アメリカ村」及び「ミナミ」が出現している。したがって、「アメリカ村」及び「ミナミ」は、発話 $U_4$ の時点で参照確率を夫々算出した場合、その値は高いはずである。「アメリカ村」も「ミナミ」も、「大阪」の代表的な繁華街であるから、発話 $U_4$ で「大阪」の単語が出現又は参照していなくとも、「アメリカ」又は「ミナミ」が出現していることによって、関連の深い「大阪」の顕現性は本来、高くなるはずである。したがって、図20の例では、発話 $U_4$ における「大阪」の顕現性を現す参照確率は、高い値を有しているはずである。

#### 【0243】

そこで、実施の形態3では、文単位又は言葉毎に対応付けられる各単語の顕現性を表わす重み値を、関連する単語(関連語)の顕現性を考慮して算出し直す。

#### 【0244】

参照確率を関連語の顕現性を考慮した重み値に算出し直すためにはまず、文単位検索装置1は、いずれの単語同士の関連が深いのかを表わす情報を先に取得しておく必要がある。そして次に、文単位毎に算出されている各単語の参照確率に、関連の深さを表わす関連度の影響を反映しておく。具体的には、例えば上述の例を用いた場合、「アメリカ村」の「大阪」への関連度を定量的に算出しておく。次に既に算出されている「アメリカ村」の参照確率へ、「大阪」への関連度の効果を反映させて、その文単位での「大阪」の顕現性を表わす重み値として算出し直して記憶しておく。

#### 【0245】

そこで、実施の形態3ではまず、文単位検索装置1は、各単語の一の単語への関連度が重み値として付与された、一の単語に対する重み付き関連語群を作成する。具体的には、実施の形態1又は2において、「3-3.文単位毎の顕現性の定量化」の処理によって文単位毎に対応付けられて記憶されている重み付き単語群、即ち単語と単語の参照確率との組又は顕現性ベクトルを利用して、文単位検索装置1が各単語の重み付き関連語群を作成する。文単位検索装置1は、文書集合全体から抽出される各単語について、夫々の単語に対する重み付き関連語群を作成し、記憶しておく。

#### 【0246】

そして次に、文単位検索装置1は、文単位毎に対応付けられて記憶されている重み付き単語群、即ち単語と単語の参照確率との組又は顕現性ベクトルの各単語の参照確率へ、各単語に関連が深い単語の参照確率からの影響を、関連度を利用して反映させ、各単語の重み値を算出し直して記憶する。

#### 【0247】

さらに、文単位検索装置1は検索処理において、各言葉に対応付けた重み付き単語群、即ち単語と単語の参照確率との組又は顕現性ベクトルについても同様に関連度を利用して

各単語の重み値を算出し直す。文単位検索装置 1 は、受け付けた言葉に対応する単語と各単語に対して算出し直した重み値に基づいて、検索処理を行なう。

【0248】

以下に、文単位検索装置 1 の CPU 11 が、各単語に対する重み付き関連語群の作成する処理について、「3-4. 関連語群の作成」の節を追加して説明する。また、作成された関連語群を使用して、「3-3. 文単位毎の顕現性の定量化」において算出した参照確率を関連を加味した重み値に算出し直す処理について、「3-5. 連想の加味した意味のまとまりの定量化」の節を追加して説明する。「4-2. 受け付けた言葉に対する意味のまとまりの定量化」において算出した参照確率を関連を加味した重み値に算出し直して検索を実行する処理について、「4-2'. 受け付けた言葉に対する連想を加味した意味のまとまりの定量化」の節を設けて説明する。

10

【0249】

なお、実施の形態 3 における、本発明に係る文単位検索装置 1 を用いた検索システムの「1. ハードウェアの構成及び概要」、及び「2. 文書データの取得及び自然言語解析」については、実施の形態 1 と同様であるため説明を省略する。「3. 文書データの文毎の意味のまとまりの定量化」、及び「4. 検索処理」について以下に説明するが、実施の形態 1 と同一の符号を用いて説明する。なお、「3. 文書データの文毎の意味のまとまりの定量化」、及び「4. 検索処理」についても、実施の形態 1 と共通する点については詳細な説明を省略する。

【0250】

#### 3-4. 関連語群の作成

関連語群は、図 6 で示した説明図で抽出されている全単語について一単語ずつ、文単位検索装置 1 によって以下の処理が行なわれることにより作成される。

20

【0251】

まず、文単位検索装置 1 は、「3-3. 文単位毎の顕現性の定量化」で全ての文単位毎に対応付けられて記憶されている重み付き単語群から、一の単語の参照確率が所定値以上の重み付き単語群を抽出する。これは、上述のように関連語を、一の単語が注目されている時に同時に注目されやすい単語とするからであり、一の単語が注目されていない文単位が除去されるようにするためである。

【0252】

次に文単位検索装置 1 は、上述の処理で抽出された、一の単語の参照確率が所定値以上の重み付き単語群を統合する。具体的には、各重み付き単語群の各単語の参照確率に、その重み付き単語群に含まれる一の単語の参照確率による重み付けをして各単語の参照確率を平均化する。一の単語の参照確率による重み付けを行うのは、一の単語の参照確率がより高い重み付き単語群の各単語に対する参照確率を使用するためである。

30

【0253】

そして、全単語についての重み付き関連語群を同様に扱うため、重み付き関連語群の各単語の重み値を正規化する。

【0254】

以下に、本発明に係る文単位検索方法を実施する文単位検索装置 1 の CPU 11 が、関連語群を作成する処理について説明する。図 2 1 及び図 2 2 は、実施の形態 3 における文単位検索装置 1 の CPU 11 が関連語群を作成する処理手順を示すフローチャートである。図 2 1 及び図 2 2 のフローチャートに示す処理は、一の単語について、その重み値が所定値以上である単語群を抽出する処理、抽出した単語群の各単語の重み値を統合して関連度として各単語に付与した関連単語群を作成する処理、一の単語に対応付けて記憶しておく処理、各単語について各処理を実行する処理に対応する。

40

【0255】

文単位検索装置 1 の CPU 11 は、記憶手段 1 3 に記憶してあるリストから一の単語を選択する(ステップ S 6 0 1)。CPU 11 は、文書記憶手段 2 から文書集合接続手段 1 6 を介してタグ付け済みの文書データを取得する(ステップ S 6 0 2)。CPU 11 は、

50

取得した文書データに付加されたタグ< s u >を文字列解析によって識別し、文単位を読み出す(ステップS 6 0 3)。次にCPU 1 1は、< s u >内に記憶してあるs a l i e n c e属性を読み出し(ステップS 6 0 4)、s a l i e n c e属性に記憶してある単語及び単語の参照確率の組(重み付き単語群)の内、ステップS 6 0 1で選択した一の単語の参照確率が所定値以上であるか否かを判断する(ステップS 6 0 5)。

【0256】

CPU 1 1が参照確率が所定値未満である(選択した一の単語が対応付けられていない)と判断した場合(S 6 0 5 : N O)、CPU 1 1は、処理をステップS 6 0 3へ戻して、後続の文単位を読み出し(S 6 0 3)、ステップS 6 0 4及びステップS 6 0 5の処理を行なう。

10

【0257】

CPU 1 1が参照確率が所定値以上であると判断した場合(S 6 0 5 : Y E S)、CPU 1 1は、ステップS 6 0 4でs a l i e n c e属性で読み出した重み付き単語群を一時記憶領域に記憶する(ステップS 6 0 6)。

【0258】

CPU 1 1は、ステップS 6 0 2で取得した文書データの全文単位についてステップS 6 0 4からステップS 6 0 6までの処理を実行したか否かを判断する(ステップS 6 0 7)。CPU 1 1が全文単位について処理を実行していないと判断した場合(S 6 0 7 : N O)、CPU 1 1は、処理をステップS 6 0 3へ戻して、後続の文単位を読み出し(S 6 0 3)、ステップS 6 0 4からステップS 6 0 6までの処理を実行する。

20

【0259】

CPU 1 1が全文単位について処理を実行したと判断した場合(S 6 0 7 : Y E S)、CPU 1 1は、全文書データについて、選択した一の単語の参照確率が所定値以上である重み付き単語群を抽出したか否かを判断する(ステップS 6 0 8)。CPU 1 1が全文書データについて選択した一の単語の参照確率が所定値以上である重み付き単語群を抽出していないと判断した場合(S 6 0 8 : N O)、CPU 1 1は、処理をステップS 6 0 2へ戻して次の文書データを取得して(S 6 0 2)ステップS 6 0 3からステップS 6 0 7までの処理を実行する。

【0260】

CPU 1 1が全文書データについて選択した一の単語の参照確率が所定値以上である重み付き単語群を抽出したと判断した場合(S 6 0 8 : Y E S)、CPU 1 1は、ステップS 6 0 6の処理によって抽出され、一時記憶領域1 4に記憶してある重み付き単語群の集合を、夫々での一の単語の参照確率で重み付けした重み値の総和を夫々の単語に対して算出することにより作成する(ステップS 6 0 9)。

30

【0261】

CPU 1 1は、ステップS 6 0 9において作成した一の単語の参照確率が所定値以上である重み付き単語群の総和、即ち総和された重み付き単語群の各単語の重み値を正規化する(ステップS 6 1 0)。

【0262】

CPU 1 1は、ステップS 6 1 0で正規化された一の単語の参照確率が所定値以上である重み付き単語群を、各重み値を関連度とする関連語群としてステップS 6 0 1で選択した一の単語に対応付けて記憶手段1 3に、又は文書集合接続手段1 6を介して文書記憶手段2に記憶する(ステップS 6 1 1)。

40

【0263】

次に文単位検索装置1のCPU 1 1は、記憶手段1 3に記憶してあるリストの全単語について関連語群を作成して記憶したか否かを判断する(ステップS 6 1 2)。CPU 1 1が全単語について関連語群を作成して記憶していないと判断した場合(S 6 1 2 : N O)、CPU 1 1は、処理をステップS 6 0 1へ戻して次の一の単語を選択し(S 6 0 1)、選択した単語についてステップS 6 0 2からステップS 6 1 1までの処理を実行する。

【0264】

50

C P U 1 1 が全単語について関連語群を作成して記憶したと判断した場合 ( S 6 1 2 : Y E S )、C P U 1 1 は処理を終了する。

【 0 2 6 5 】

なお、ステップ S 6 0 5 において文単位検索装置 1 の C P U 1 1 は、単純に、参照確率が所定値以上であるか否かを判断するのではなく、以下のような正規化処理を行ってから所定値との比較を行うようにしてもよい。例えば、文単位検索装置 1 の C P U 1 1 は、文単位に対応付けられている各単語の参照確率の二乗の総和が「1」になるように、全参照確率の二乗和の二乗根で各参照確率を除算することによって正規化を行う。

【 0 2 6 6 】

なお、ステップ S 6 1 0 における正規化についても、各単語の重み値の二乗の総和が 1 になるように正規化する。例えば、文単位検索装置 1 の C P U 1 1 は、全重み値の二乗和の二乗根により、各重み値を除算することによって正規化を行う。

【 0 2 6 7 】

次に、実施の形態 3 における文単位検索装置 1 の C P U 1 1 が、図 2 1 及び図 2 2 のフローチャートに示した処理を一の単語について行った場合に作成される関連語群の具体例を示す。

【 0 2 6 8 】

図 2 3 は、実施の形態 3 における文単位検索装置 1 の C P U 1 1 によって関連語群が作成される場合の、各処理の過程での重み付き単語群の例を示す説明図である。なお、図 2 3 の説明図に示す例は、文単位検索装置 1 の C P U 1 1 によって、一の単語「アメリカ村」の参照確率が所定値 ( 0 . 2 ) 以上の重み付き単語群が抽出された場合の例である。図 2 3 ( a ) は、図 2 1 及び図 2 2 のフローチャートに示したステップ S 6 0 5 における C P U 1 1 の処理により抽出されて、一時記憶領域 1 4 に記憶されている重み付き単語群  $G W_1$  ,  $G W_2$  ,  $G W_3$  を示している。図 2 3 ( b ) は、同様にステップ S 6 0 7 における C P U 1 1 の処理により、一の単語の参照確率により重み付けされる重み付き単語群  $G W_1'$  ,  $G W_2'$  ,  $G W_3'$  を示している。図 2 3 ( c ) は、同様にステップ S 6 0 9 における C P U 1 1 の処理により、重み付けされて総和された重み付き単語群  $G W''$  を示している。

【 0 2 6 9 】

図 2 3 ( a ) に示すように、一の単語「アメリカ村」の重み値 ( 参照確率 ) が所定値 0 . 2 以上の重み付き単語群  $G W_1$  ,  $G W_2$  ,  $G W_3$  が抽出されている。

【 0 2 7 0 】

図 2 3 ( b ) に示されている重み付き単語群  $G W_1'$  ,  $G W_2'$  ,  $G W_3'$  の、各単語の重み値には夫々の重み付き単語群中の一の単語「アメリカ村」の重み値 ( 参照確率 ) が乗算されている。図 2 3 ( a ) に示された単語群  $G W_1$  ,  $G W_2$  ,  $G W_3$  に対し、図 2 3 ( b ) に示された単語群  $G W_1'$  ,  $G W_2'$  ,  $G W_3'$  の各単語の重み値は、以下のようにして一の単語「アメリカ村」の重み値 ( 参照確率 ) が乗算されている。例えば、重み付き単語群  $G W_1$  の各単語の重み値は、アメリカ村の重み値 ( 参照確率 ) が 0 . 6 であるため、アメリカ村の参照確率で重み付けされて以下のようになる。

【 0 2 7 1 】

単語群  $G W_1'$  : ( 秋 : 0 ( 0 . 6 × 0 ) , アメリカ村 : 0 . 3 6 ( 0 . 6 × 0 . 6 ) , . . . , 大熊座 : 0 ( 0 . 6 × 0 ) , 大阪 : 0 . 1 2 ( 0 . 6 × 0 . 2 ) , 大鹿 : 0 ( 0 . 6 × 0 ) , . . . )

【 0 2 7 2 】

つまり、一の単語「アメリカ村」の重み値が高いほど、他の単語の重み値の影響が反映される。

【 0 2 7 3 】

図 2 3 ( c ) に示されている重み付き単語群  $G W''$  の、各単語の重み値は、図 2 3 ( b ) に示したように夫々一の単語「アメリカ村」の重み値 ( 参照確率 ) で重み付けされた重み値が単語毎に総和されている。図 2 3 ( c ) に示された単語群  $G W''$  の各単語の重

10

20

30

40

50

み値は、図 23 (b) に示された単語群  $GW_1'$ ,  $GW_2'$ ,  $GW_3'$  以下のように総和される。

【0274】

単語群  $GW'$  : (秋: 0.03 (= 0 + 0.03 + 0), アメリカ村: 0.49 (= 0.36 + 0.09 + 0.04), ..., 大熊座: 0 (= 0 + 0 + 0), 大阪: 0.28 (= 0.12 + 0.12 + 0.04), 大鹿: 0 (= 0 + 0 + 0), ...)

【0275】

また、重み付けされて総和されることにより統合された重み付き単語群  $GW'$  の各単語の重み値は、文単位検索装置 1 の CPU 11 の処理により正規化される。

【0276】

正規化の処理についてはその方法は問わないが、例えば、文単位検索装置 1 の CPU 11 は、各単語の重み値を二乗し、二乗した値の和の二乗根を算出し、各単語の重み値で割って、重み付き単語群  $GW'$  の各単語の重み値を正規化するようにしてもよい。

【0277】

また、重み付けされて総和されることにより統合された重み付き単語群  $GW'$  を、各単語を一次元とし、各単語の重み値を各次元方向の要素として多次元ベクトルである関連度ベクトルで表現した場合は、各重み値(要素)を多次元ベクトルのノルムで割ることにより、多次元ベクトルを正規化するようにしてもよい。このとき、ノルムはユークリッドノルムとは限らない。

【0278】

このように総和して正規化した結果の重み付き単語群が、文単位検索装置 1 の CPU 11 により「アメリカ村」の関連語群として作成される。以下に示す例は、単語「アメリカ村」の関連語群の一例である。なお、各単語は、重み値の大きい順に列挙されている。

【0279】

関連語群(「アメリカ村」) = (アメリカ村: 0.647, アメリカ: 0.369, 大阪: 0.258, 村: 0.159, 防犯カメラ: 0.139, カメラ: 0.139, チェックアウト: 0.129, アウト: 0.129, 中: 0.128, 女性: 0.120, 男: 0.102, 中央: 0.098, 犯行: 0.092, 人: 0.087, たこ焼き: 0.082, 心齋橋: 0.075, ミナミ: 0.074, 警察: 0.073, 時間: 0.071, 公園: 0.065, 昭和: 0.064, 今回: 0.063, 数: 0.061, なんば: 0.060, 御津: 0.060, ランドローバー(登録商標): 0.059, ローバー(登録商標): 0.059, 名前: 0.059, プラン: 0.057, 道頓堀: 0.055, 立川: 0.055, ナンバー: 0.054, 西鉄: 0.053, サツ: 0.052, 伊那: 0.050, オリジナルステッカー: 0.049, ステッカー: 0.049, イン心齋橋: 0.049, 御堂筋線: 0.049, ...)

【0280】

なお上の例は、文書集合(GDA タグ付き毎日新聞コーパス <http://www.gsk.or.jp/catalog.html> 参照)を使用して実際に作成した「アメリカ村」の関連語群である。

【0281】

上述の「アメリカ村」の関連語群の具体例に示したように、例えば、「アメリカ村」が注目されている場合、「大阪」は他の単語よりも注目される関連語であることを重み値によって定量的に表わすことができる。したがって、この関連語群の各単語の重み値は一の単語への関連度を表わしているといえることができる。上述の具体例では「アメリカ村」の「大阪」への関連度は、0.258 である。

【0282】

以下、単語  $w_j$  に対して作成した関連語群の各重み値、即ち単語  $w_j$  の単語  $w_k$  への関連度を  $b_{j,k}$  と表わす。一の単語  $w_j$  の関連語群は  $bw_j = (w_1 : b_{j,1}, w_2 : b_{j,2}, \dots, w_n : b_{j,n})$  と表わされる。なお、関連語群を関連度ベクトルとして表わす場合、 $bw_j = (b_{j,1}, b_{j,2}, \dots, b_{j,n})$  と表現する。

【0283】

10

20

30

40

50



文単位検索装置1のCPU11は、上述のような処理を、図6の説明図に示した全単語について繰り返し行って各単語の関連単語群を作成し、文書記憶手段2又は文単位検索装置1の記憶手段13に記憶しておく。このように、文書集合に出現する単語全てについて夫々関連度が定量的に算出されて付与された関連語群を作成して記憶しておくことにより、文単位毎の意味のまとまりを表わす重み付き単語群に対し、関連語の関連度による影響を反映させることができる。

【0284】

3-5. 連想を加味した意味のまとまりの定量化

次に、文単位毎に記憶されている重み付き単語群、即ち単語と各単語の参照確率との組又は顕現性ベクトルに、作成された関連語群の各単語の関連度を反映させる。具体的には、文単位検索装置1は、既に算出されて記憶されている各単語の参照確率を読み出し、一の単語の重み値として、各単語の参照確率に各単語から一の単語への関連度を乗算した値を算出し直して記憶する。

10

【0285】

図24は、実施の形態3における文単位検索装置1のCPU11が、各文単位に対応付けられて記憶されている重み付き単語群の各単語の重み値を算出し直す処理手順を示すフローチャートである。図24のフローチャートに示す処理は、文単位毎に対応付けられた重み付き単語群の各単語の重み値を、関連度を使用して付与し直す処理に対応する。

【0286】

文単位検索装置1のCPU11は、文書記憶手段2から文書集合接続手段16を介してタグ付け済みの文書データを取得する(ステップS71)。CPU11は、取得した文書データに付加されたタグ<su>を文字列解析によって識別し、文単位を読み出す(ステップS72)。

20

【0287】

次にCPU11は、<su>内に記憶してあるsalience属性を読み出し(ステップS73)、salience属性で対応付けて記憶してある単語及び単語の参照確率の組(重み付き単語群)の、各参照確率を関連語群を使用して連想を加味した重み値に算出し直す(ステップS74)。CPU11は、各単語及び各単語についてステップS74で算出し直した重み値の組である重み付き単語群(顕現性ベクトル)をsalience属性を付加して記憶し直す(ステップS75)。

30

【0288】

次にCPU11は、ステップS72で読み出した文単位が文書データの終端であるか否かを判断する(ステップS76)。現在の文が取得した文書データの終端であるか否かは、現在の文を挟む<su></su>の後に、<su>タグが後続するかしないかを判断し、後続しないと判断した場合は終端であると判断することができる。CPU11が文書データの終端でないと判断した場合は(S76:NO)、CPU11は、処理をステップS72に戻し、次の文単位に対して処理を継続する。一方、CPU11が文書データの終端であると判断した場合は(S76:YES)、CPU11は、全文書データについて、重み付き単語群の各単語の重み値を算出し直してsalience属性で対応付けて記憶する処理を終了したか否かを判断する(ステップS77)。

40

【0289】

CPU11が全文書データについて、重み付き単語群の各単語の重み値を算出し直してsalience属性によって記憶する処理を終了していないと判断した場合は(S77:NO)、CPU11は、処理をステップS71へ戻し、別の文書データを取得して処理を継続する。CPU11が全文書データについて、重み付き単語群の各単語の重み値を算出し直してsalience属性によって記憶する処理を終了したと判断した場合は(S77:YES)、CPU11は処理を終了する。

【0290】

なお、文単位検索装置1のCPU11は、ステップS74における各単語の重み値の算出し直しを以下のような処理を行なうことによって実現する。

50

## 【0291】

図25は、実施の形態3における文単位検索装置1のCPU11が、各文単位に対応付けられて記憶されている重み付き単語群の各単語の重み値を算出し直す処理手順の詳細を示すフローチャートである。図25のフローチャートに示す処理は、各単語の関連度を重み付き単語群の重み値に乘算する処理、乗算した重み値に基づいて各単語の重み値を付与し直す処理に対応する。

## 【0292】

文単位検索装置1のCPU11は、図24のフローチャートのステップS74で読み出したsalience属性で対応付けて記憶してある重み付き単語群の各単語及び各単語の参照確率を読み出し、一時記憶領域14に記憶しておく(ステップS81)。CPU11は、各単語の内の一の単語を選択し(ステップS82)、選択した一の単語の重み値について以下の処理を行なう。

## 【0293】

CPU11は、記憶手段13又は文書記憶手段2に記憶してある各単語の関連度が付与された関連語群を読み出す(ステップS83)。CPU11は、読み出した各単語の関連語群から、各単語から一の単語への関連度を取得する(ステップS84)。CPU11は、取得した各単語から一の単語への関連度を一時記憶領域14に記憶してある各単語の参照確率に夫々乗算し、和を算出する(ステップS85)。

## 【0294】

CPU11によりステップS85で算出された和が、一の単語について、関連語による連想が加味されて算出し直された顕現性を表わす重み値である。

## 【0295】

CPU11は、ステップS81で一時記憶領域14に記憶してある各単語全てについて、重み値を算出し直したか否かを判断する(ステップS86)。CPU11が各単語全てについて重み値を算出し直していないと判断した場合(S86:NO)、CPU11は、処理をステップS82へ戻して、次の単語についてステップS82からステップS85までの重み値を算出し直す処理を実行する。CPU11が各単語全てについて重み値を算出し直したと判断した場合(S86:YES)、CPU11は、処理を図24のフローチャートのステップS75へ戻す。

## 【0296】

なお、図24のフローチャートの内のステップS74及び図25のフローチャートに示した文単位検索装置1のCPU11による重み値を算出し直す処理は、実施の形態1における参照確率を算出して各文単位毎の顕現性を現す重み値として記憶する処理の中で実行してもよい。具体的には、図9のフローチャートに示した処理手順の内のステップS306とステップS307の処理の間にステップS74及び図25のフローチャートに示した処理を実行する構成でもよい。

## 【0297】

図24及び図25のフローチャートに示したCPU11の処理手順において、文単位検索装置1のCPU11が、各単語について算出した参照確率を連想を加味した重み値に算出し直す処理について、具体的な例を以下に示す。

## 【0298】

例えば、単語「アメリカ村」について作成した関連度群を使用する場合、文単位検索装置1により、ある文単位における「大阪」の顕現性を現す重み値を以下のように算出し直す。なお、「アメリカ村」について作成した関連度群の「大阪」への関連度は「0.3」であるとする。ある文単位に対応付けて記憶されている単語に「アメリカ村」が含まれており、「アメリカ村」の参照確率が0.4であり、「大阪」は含まれていない場合であっても、文単位検索装置1のCPU11は、「アメリカ村」の参照確率0.4に、「アメリカ村」から「大阪」への関連度0.3を乗算して、その文単位における「大阪」の重み値は「0」ではなく「0.12」に算出し直す。

## 【0299】

10

20

30

40

50

ここで、文脈連想を加味した単語  $w_k$  の各文  $s_i$  における顕現性を表わす重み値を、 $salience(w_k | pre(s_i))$  と表わす。また、単語  $w_k$  の各文  $s_i$  における参照確率を  $Pr(w_k | pre(s_i))$  とする。この場合、単語  $w_j$  の単語  $w_k$  への関連度を反映した場合、 $salience(w_k | pre(s_i)) = b_{j,k} \times Pr(w_j | pre(s_i))$  と算出し直される。なお、単語  $w_k$  への関連度を有する単語  $w_j$  は他にも存在するので、全単語  $w_j$  ( $j = 1, \dots, N$ ) からの関連度の影響をも反映させて、文単位検索装置 1 は以下に示す式 (6) のように各単語の重み値を算出し直す。

【0300】

【数6】

$$salience(w_k | pre(s_i)) = \sum_{j=1}^N b_{j,k} \times Pr(w_j | pre(s_i)) \quad \dots (6)$$

10

【0301】

したがって、文単位検索装置 1 の CPU 11 は、以下に示す式 (7) のように文単位  $S$  における各単語  $w_k$  ( $k = 1, \dots, N$ ) の重み値を算出し直す。

20

【0302】

【数7】

$$\begin{aligned} \vec{V}(S_i) &= (salience(w_1 | pre(s_i)), \dots, salience(w_k | pre(s_i)), \dots, salience(w_N | pre(s_i)))^T \\ &= \begin{bmatrix} b_{1,1} & \dots & b_{N,1} \\ \vdots & b_{j,k} & \vdots \\ b_{1,N} & \dots & b_{N,N} \end{bmatrix} \begin{bmatrix} Pr(w_1 | pre(s_i)) \\ \vdots \\ Pr(w_j | pre(s_i)) \\ \vdots \\ Pr(w_N | pre(s_i)) \end{bmatrix} \\ &= (\vec{b}_{w_1}, \dots, \vec{b}_{w_k}, \dots, \vec{b}_{w_N}) \vec{v}(s_i) \quad \dots (7) \end{aligned}$$

30

【0303】

なお、式 (7) の最終行の式は、実施の形態 2 に示したように、重み付き単語群、即ち単語と単語の参照確率との組を顕現性ベクトル  $v(s_i)$  として表現した場合に、 $salience(w_k | pre(s_i))$  を  $k$  番目の要素として有する連想を加味した後の顕現性ベクトル  $V(s_i)$  の各単語の重み値の算出の原理を表わす。

40

【0304】

この場合、各  $b_{w_1}, \dots, b_{w_N}$  は、全単語  $w_1, \dots, w_N$  に対する関連語群をベクトルによって表現した関連度ベクトルである。

【0305】

重み付き単語群、即ち単語と単語の参照確率との組を多次元ベクトル  $v(s_i)$  で表現し、関連語群を関連度ベクトル  $b_{w_1}, \dots, b_{w_N}$  で表現した場合、式 (7) のように各単語の参照確率を、連想を加味した重み値に算出し直す処理は、以下のように解釈するこ

50

とができる。

【0306】

$salience(w_k | pre(s_i))$  を  $k$  番目の要素として有する、連想を加味した顕現性ベクトル  $V(s_i)$  は、関連度ベクトル  $bw_1, \dots, bw_N$  を基底とする斜交座標系における顕現性ベクトル  $v(s_i)$  であると解釈することができる。言い換えると、連想を加味した顕現性ベクトル  $V(s_i)$  は、参照確率をそのまま要素とする顕現性ベクトル  $v(s_i)$  を関連語軸方向へ回転させたものであると解釈することができる。

【0307】

関連度ベクトル  $bw_1, \dots, bw_N$  を基底とする斜交座標系とは、連想を加味した各単語を1次元とした場合に、各基底ベクトル(各単語の次元方向に大きさ1のベクトル)は、夫々直行せず関連度が高い単語同士の基底ベクトル間の角度が小さくなるような座標系である。

10

【0308】

$b_{j,k}$  を各要素とする変換行列を参照確率を要素とする顕現性ベクトルに乗算すると、関連する単語の次元方向に回転した顕現性ベクトル  $V(s_i)$  が得られると解釈することができる。

【0309】

したがって、文毎の意味のまとまりを表わす重み付き単語群を顕現性ベクトルで表現して記憶している場合、文単位検索装置1のCPU11がその顕現性ベクトルを関連度ベクトルによって回転(変換)する処理を行なうことによって、文毎の意味のまとまりを連想が加味された顕現性ベクトルで表わして記憶しておくことができる。

20

【0310】

次に、上述のように定量的に関連度を表わした関連度群を使用して、各文単位の意味のまとまりを表わす各単語の重み値を連想を加味して算出し直す処理を実行した結果の具体例を以下に示す。図26は、実施の形態3における文単位検索装置1のCPU11によって算出された各単語の顕現性を表わす重み値の内容例を示す説明図である。図26(a)に示した各文  $s_1, s_2$  に対する各単語の重み値は夫々、関連語群を使用して連想が加味される前の参照確率の値である。一方、図26(b)に示した各文  $s_1, s_2$  に対する各単語の重み値は、関連語群を使用して連想が加味された後の重み値である。

【0311】

30

なお、図26に示す具体例は、日本語話し言葉コーパス([http://www.kokken.go.jp/katsudo/kenkyu\\_jyo/corpus/](http://www.kokken.go.jp/katsudo/kenkyu_jyo/corpus/)、CSJ/vol17/D03F0040)より抽出した文単位の例である。

【0312】

図26の内容例に示すように、図26(b)の文  $s_1$  における「大阪」の重み値は、図26(a)の文  $s_1$  における「大阪」の参照確率の値  $0.3338$  と比較して、 $0.6229$  と高くなっている。また、図26(b)の文  $s_2$  における「大阪」の重み値は、図26(a)の文  $s_2$  における参照確率の値  $0.3208$  と比較して、 $0.6675$  とさらに高くなっている。

【0313】

40

さらに、図26(a)の参照確率の例では、文  $s_2$  における「大阪」の重み値は、文  $s_2$  に「アメリカ村」が出現しているにも拘わらず、その「大阪」の重み値への影響(励起)が考慮されていないために重み値が低下している。これに対し、図26(b)の連想を加味した後の重み値の例では、文  $s_2$  における「大阪」の重み値は、文  $s_2$  に「アメリカ村」が出現していることによって、出現していない「大阪」の顕現性を表わす重み値が高くなっている。「アメリカ村」と「大阪」との関連度の影響が反映されているからである。

【0314】

このように、文単位検索装置1が文単位毎に記憶している重み付き単語群に対し、参照確率という定量的な値を用いて関連度を表わした関連語群を用いて連想を加味することに

50

より、文単位で「アメリカ村」が注目されている場合の「大阪」の顕現性を、文単位又は言葉の書き手又は話し手の背景文脈により近づかせることができる。これにより、「大阪」の単語の顕現性を表わす重み値が低く算出されて、文単位の意味のまとまりが書き手又は話し手の実際の文脈と離れたように定量的に評価されてしまうことを回避することができる。

#### 【0315】

##### 4. 検索処理

次に、実施の形態3における検索処理について説明する。「4-1. ユーザから入力された言葉の受け付け」については、受付装置4のCPU41が行う処理については実施の形態1及び2と同様であるので、詳細な説明を省略する。

10

#### 【0316】

##### 4-2'. 受け付けた言葉に対する連想を加味した意味のまとまりの定量化

次に、文単位検索装置1のCPU11が、受付装置4, 4, ...で受け付けた言葉のデータを受信した場合に、文書記憶手段2で記憶している文書中の文を検索する処理について説明する。受け付けた言葉に対しても、意味のまとまりの定量化、即ち当該テキストデータの単語抽出及び単語の参照確率を算出し、さらに関連度を使用して重み値を算出し直す。

#### 【0317】

実施の形態3では、文単位検索装置1のCPU11は、受け付けた言葉の意味のまとまりを定量的に表わす単語と単語の参照確率との組又は顕現性ベクトル、即ち重み付き単語群に、関連語による連想を加味する。以下に、文単位検索装置1のCPU11が受け付けた言葉に対応付けた重み付き単語群の各単語の重み値を連想を加味して算出し直し、算出し直した重み値に基づいて検索を実行する処理について説明する。

20

#### 【0318】

図27は、実施の形態3における文単位検索装置1及び受付装置4の検索処理の処理手順を示すフローチャートである。なお、図27のフローチャートに示す処理手順では、実施の形態1における図15、図16及び図17のフローチャートに示した検索処理の処理手順と同一の処理については各ステップに同一の符号を用いて詳細な説明を省略する。

#### 【0319】

図27のフローチャートに示す処理手順の内、二点鎖線で囲まれたステップS4001の処理が、実施の形態1における図15、図16及び図17のフローチャートに示した処理手順と異なる。即ち、ステップS411と、ステップS412との間に以下に説明するステップS4001が追加されていることが異なる。

30

#### 【0320】

以下に、実施の形態3において受け付けた言葉の意味のまとまりを表わす重み付き単語群を対応付け、予め記憶してある意味のまとまりが類似する文単位を抽出する検索処理について以下に説明する。

#### 【0321】

CPU11は、一時記憶領域14に夫々参照確率を算出して記憶している全単語に対し、所定値以上の参照確率が算出された単語に絞り込み(ステップS411)、ステップS408において算出した参照確率を、連想を加味した重み値に算出し直す(ステップS4001)。ステップS4001における、CPU11による連想を加味した重み値の算出し直しの処理は、図25のフローチャートに示した処理と同様、単語を1つずつ選択し、選択した一の単語への各単語の関連度と各単語の参照確率とを乗算して算出する。

40

#### 【0322】

それまでの処理により、受け付けた言葉に対し、以前に受け付けた言葉から続く流れ上の意味のまとまりを、連想を加味した上で定量的に表わす単語と単語の参照確率の組(重み付き単語群)を検索要求として生成することができた。

#### 【0323】

CPU11はこの後、ステップS4001で得られた連想が加味された重み付き単語群

50

に対し、各文毎に対応付けて記憶してある、連想が加味された重み付き単語群を読み出して、類似する文を抽出する処理を実行する。連想が加味された重み付き単語群についての以降の処理は実施の形態 1 と同様であるので詳細な説明を省略する。

【0324】

これにより、文単位検索装置 1 は、文書記憶手段 2 に記憶してある文書データから分別される文と受け付けた言葉とで、関連語を利用して連想を加味した意味のまとまりが類似しているか否かを判断し、類似すると判断された文を直接的に出力することができる。したがって、本発明の文単位検索方法を実施することにより、文脈上の意味のまとまりが類似する文単位を連想を加味して効果的に抽出し、直接的に出力することができる。

【0325】

なお、文単位検索装置 1 の CPU 11 は、受け付けた言葉に対して重み付き単語群を対応付け、文毎に予め記憶してある重み付き単語群と類似しているか否かを判断する場合、図 27 のフローチャートに示した処理手順のように、重み付き単語群が同一の単語を含んでいるか否かによって判断するとは限らない。さらに同一の単語に付与されている重み値の差分を算出し、算出した差分が小さい程類似すると判断するとは限らない。

【0326】

次に、文単位検索装置 1 の CPU 11 が、受け付けた言葉と意味のまとまりが類似する文単位を抽出する処理を、意味のまとまりを顕現性ベクトル及び関連度ベクトルで表現し、ベクトル間の距離を算出することによって実現する場合について以下に説明する。

【0327】

図 28 は、実施の形態 3 におけるベクトル表現を用いた場合の文単位検索装置 1 及び受付装置 4 の検索処理の処理手順を示すフローチャートである。なお、図 28 のフローチャートに示す処理手順では、実施の形態 1 における図 15、図 16 及び図 17 のフローチャート、及び実施の形態 2 における図 19 のフローチャートに示した検索処理の処理手順と同一の処理については各ステップに同一の符号を用いて詳細な説明を省略する。

【0328】

図 28 のフローチャートに示す処理手順の内、一点鎖線で囲まれた各ステップ S501 からステップ S506 までの処理が、実施の形態 1 における図 15、図 16 及び図 17 のフローチャートに示した処理手順と異なる。実施の形態 1 におけるステップ S412 からステップ S416 までの処理の代わりに、実施の形態 2 における文単位検索装置 1 の CPU 11 により実行されるステップ S501 からステップ S506 までの処理と同様の処理を行なう。図 28 のフローチャートに示す処理手順の内、二点鎖線で囲まれたステップ S5001 の処理が、実施の形態 2 における図 19 のフローチャートに示した処理手順と異なる。即ち、ステップ S501 と、ステップ S502 との間に以下に説明するステップ S5001 が追加されていることが異なる。

【0329】

文単位検索装置 1 の CPU 11 は、ステップ S501 で算出した顕現性ベクトルを、関連語による連想を加味した顕現性ベクトルに算出し直す(ステップ S5001)。

【0330】

CPU 11 はこの後、ステップ S5001 で得られた連想が加味された重み付き単語群に対し、各文毎に対応付けて記憶してある、連想が加味された顕現性ベクトルを読み出して、類似する文を抽出する処理を実行する。連想が加味された顕現性ベクトルを読み出して類似する文を抽出する処理は実施の形態 2 と同様であるので詳細な説明を省略する。

【0331】

なお、CPU 11 によるステップ S5001 において、顕現性ベクトルを関連語による連想を加味した顕現性ベクトルに算出し直す処理は、ステップ S501 で算出した顕現性ベクトルを関連度ベクトル群(行列)で式(7)で示したように変換して(回転させて)算出する。具体的には、参照確率のみを要素とする多次元ベクトル  $v(s_i)$  に対して上述の連想を加味した顕現性ベクトル  $V(s_i)$  を算出する。

【0332】

10

20

30

40

50

なお、上述の図 2 8 のフローチャートに示した処理手順の内の、CPU 11 が受け付けた言葉に対応付けた顕現性ベクトルと、読み出した顕現性ベクトルとの距離を算出するステップ S 5 0 3 の処理は、実施の形態 3 では、具体的には以下のように算出する。受け付けた言葉  $u_i$  に対し連想が加味されて算出し直された顕現性ベクトルが  $V(u_i)$  と表わされ、読み出された、予め連想が加味されてある顕現性ベクトルが  $V(s_i)$  と表わされる場合、CPU 11 は以下に示す式 (8) のように、コサイン距離を算出する。

【0333】

【数 8】

$$\frac{\vec{V}(s_i) \cdot \vec{V}(u_i)}{|\vec{V}(s_i)| |\vec{V}(u_i)|}$$

$$= \frac{\sum_{k=1}^N \text{saliency}(w_k | s_i) \text{saliency}(w_k | u_i)}{\sqrt{\sum_{k=1}^N \text{saliency}(w_k | s_i)^2} \sqrt{\sum_{k=1}^N \text{saliency}(w_k | u_i)^2}} \dots (8)$$

10

20

【0334】

ただし、式 (8) に示したように距離を算出した場合、言葉の顕現性ベクトル  $V(u_i)$  と、読み出した顕現性ベクトル  $V(s_i)$  とが近いほど、算出したコサイン距離の値は大きくなる。したがって、CPU 11 はステップ S 5 0 6 において、算出したコサイン距離が大きい順に類似度を付与する。

【0335】

文単位検索装置 1 の CPU 11 による上述のような処理により、連想が加味された意味のまとまりを表わす顕現性ベクトル間の距離によって、意味のまとまりが類似する文単位を直接的に検索することができる。ベクトル表現を用いることにより、CPU 11 は、受け付けた言葉に対応付けられる連想が加味された重み付き単語群と、予め文に対応付けて記憶されている連想が加味された重み付き単語群とを一単語ずつ重み値を比較している処理を行なうことなしに、連想を加味した上で直接的に類似しているか否かを判断を行うことができる。

【0336】

また、実施の形態 3 における文単位検索装置 1 による場合、各文単位及び単語に対応付けられる顕現性ベクトルは、各単語に相当する次元間が直交しない関連度が高い単語の次元方向間の角度が小さくなるような斜交座標系で扱われる。このため、類似するか否かを判断する際にベクトル間の距離を比較した場合に、関連度が高い単語の次元方向に要素を有している場合は類似していると判断されるようになる。

【0337】

したがって、「大阪」の顕現性が高い文単位  $s$  が記憶されている場合、受け付けた言葉において例えば「オランダ村」の顕現性が高いときは、文単位  $s$  は受け付けた言葉に類似すると判断されない。しかし、受け付けた言葉において「アメリカ村」の顕現性が高いときは、受け付けた言葉において「大阪」の顕現性が励起されて高くなるので、文単位  $s$  はこの受け付けた言葉に類似すると判断される可能性が高くなる。

【0338】

30

40

50

これにより、受け付けた言葉に対し、連想を加味してより効果的に意味のまとまりが類似する文単位を検索して直接的に出力することができる。

【0339】

なお、実施の形態1乃至3では、検索結果として受信したテキストデータは、受付装置4が備える表示手段46のモニタ等で表示する構成としたが、受信したテキストデータから音声に変換して、音声入出力手段47のスピーカ等を介して出力する構成でもよい。これにより、ユーザは自分が音声入力した複数の言葉によって、又は他のユーザとの会話を音声入力することで、その会話の文脈と意味のまとまりが類似する文を検索結果として得ることができる。受け付けた言葉が話し言葉からなる場合に、発話では省略されている、ゼロ代名詞で表される単語をも含めた単語の顕現性が類似する文を直接的に検索結果として得ることができる。

10

【0340】

また、文単位検索装置1のCPU11は、言葉のテキストデータを受信する都度、当該テキストデータに対して検索された文のうち、一番優先順位の高い文を表すテキストデータのみを受付装置4, 4, ...に送信する構成としてもよい。これにより、入力される言葉に対する検索結果を会話の第三者の発話として提示し、鼎談を実現することも可能である。

【0341】

なお、実施の形態1乃至実施の形態3では、文単位検索装置1は文毎に顕現性を示す情報を特定して記憶したが複数の文からなる段落(paragraph)毎にタグ<p></p>で挟み、当該段落に対して特徴パターンを特定して顕現性を示す情報をsalience属性によって記憶させ、段落を検索結果として出力する構成としてもよい。文又は段落に限らず、一定の意味のまとまりを表す単位であれば文節であっても構わない。話し言葉の場合は一文と識別できる文字列が非常に長くなることが考えられる。多数の文節から構成され、文節と文節は「～も」「～ので」等の接続助詞で続いているにも拘わらず、文脈が動的に変化していく場合は一文では意味がまとまっていないときがある。したがって、所定の文節の数を超えて構成される文の場合は、文節毎に一文であるとみなして処理を行う構成としてもよい。

20

【0342】

また、実施の形態1乃至実施の形態3では、話し言葉からなる文書データを書き言葉からなる文書データと区別して予め記憶しておく構成としたが、受信した言葉に対して各単語の特徴パターンを特定して参照確率を算出する都度、文書記憶手段2で記憶する構成としてもよい。この際、文単位検索装置1のCPU11は、連続して受信した言葉が一連のものであるか否かの判断を当該言葉の送信元である受付装置4を識別する情報と、受付装置4がユーザの検索開始・完了操作を検知したことを示す情報とによってすることもできる。これにより、予め文書記憶手段2で記憶してある文書データのページに該当する単位で言葉を文書記憶手段2に記憶させることができる。

30

【0343】

なお、実施の形態1乃至実施の形態3では、文書データの取得とタグ付け、参照確率を求めるための回帰分析、更に言葉を受け付けた際の処理を文単位検索装置1が全て行う構成としたが、文単位検索装置と文書記憶装置とに分ける構成としてもよい。この場合は、文書記憶装置でWebクロウリングを行って文書データを取得し、さらに形態素解析及び統語解析によってテキストデータにタグを付加して記憶しておく。また、文書記憶装置で記憶した文書データをもとに参照確率を算出するための式を回帰分析によって求め、求めた式を使用して、記憶した文書データに対して文毎の単語及び単語の参照確率を記憶する処理を予めしておく。文単位検索装置は、言葉を変換したテキストデータを受信した際に特徴パターンを特定し、文書記憶装置から参照確率を算出するための回帰式を取得して参照確率を算出して検索を行う。

40

【0344】

また、実施の形態1乃至実施の形態3では、ユーザからの文字列入力又は音声入力等の

50



言葉の入力は、受付装置 4 によってテキストデータに変換され、文単位検索装置 1 へ送信される構成とした。これに限らず、文単位検索装置 1 が、ユーザの文字列入力操作を受け付ける入出力手段、及びユーザの音声入力を受け付ける音声入力手段を備える構成でもよい。図 29 は、本発明の文単位検索方法を文単位検索装置 1 で実施する場合の構成を示すブロック図である。この場合、文単位検索装置 1 は、CPU 11、内部バス 12、記憶手段 13、一時記憶領域 14、文書集合接続手段 16 及び補助記憶手段 17 の他に、ユーザの操作を受け付けるマウス又はキーボード等の操作手段 145、モニタ等の表示手段 146 及びマイク及びスピーカ等の音声入出力手段 147 を更に備える。

【0345】

図 29 の構成図に示した構成の場合、文単位検索装置 1 の CPU 11 は、音声入力手段から入力された音声の特徴を表わす、周波数又は会話速度等を検知し、発話における各単語の特徴パターンを特定することができる。各単語の文法的な特徴パターンは、入力された音声を音声認識によりテキストデータに変換して当該テキストデータに基づいて検索する構成としてもよい。

10

【0346】

実施の形態 1 乃至実施の形態 3 では、受付装置 4, 4, ... は、受け付けた文字列又は音声の言葉を一定の長さに区切ってデジタルデータに変換して送信するのみの装置として構成した。しかしながら、本発明の文単位検索方法を実施するためには、受付装置 4, 4, ... が記憶手段 43 に記憶しているプログラムを、受付装置 4, 4, ... が受け付けた言葉を形態素解析及び統語解析、又は音素解析等の自然言語解析を実行することができるように構成してもよい。この場合、受付装置 4, 4, ... の CPU 41 は、受け付けた言葉における各単語の顕現性を表わす重み値を算出し、算出した重み付き単語群を検索要求として文単位検索装置 1 へ送信する構成でもよい。

20

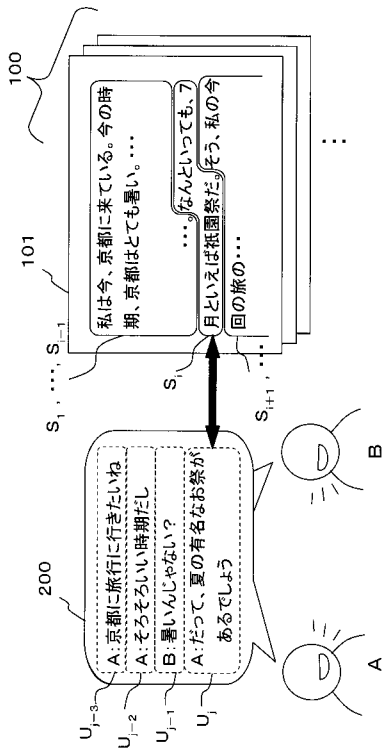
【産業上の利用可能性】

【0347】

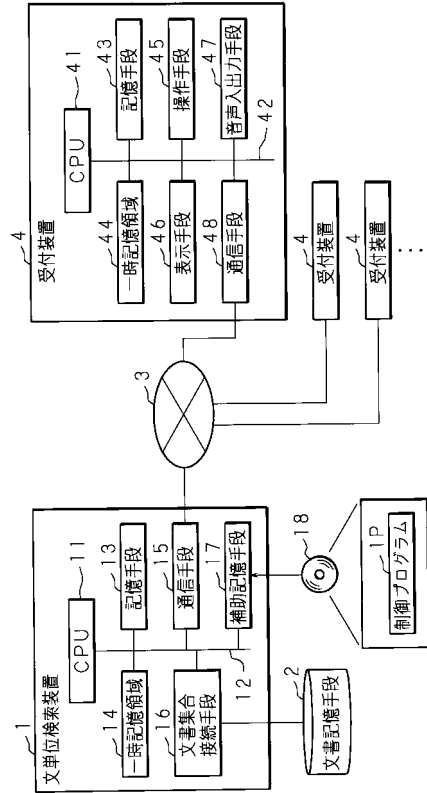
本発明に係る文単位検索方法を、ユーザ間の会話を音声認識が可能なコンピュータ装置に実施させることにより、コンピュータ装置にユーザ間の会話に参加させて鼎談を実現する用途にも適用することが可能である。また、ユーザ間の会話又はチャットの文脈の流れに応じて切り替わる会話連動型広告の提示サービスを実現する用途にも適用可能である。会議中の文脈の流れに応じて、過去の議事録から類似関連する議事録を提示する会議支援サービスへの適用も可能である。さらに、執筆中の文章を言葉として受け付け、文脈の流れに応じて、関連する情報を提供する文章執筆支援サービスへの適用も可能である。

30

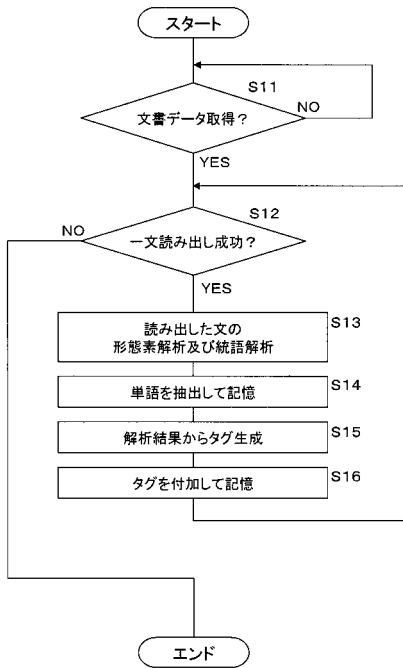
【図1】



【図2】



【図3】



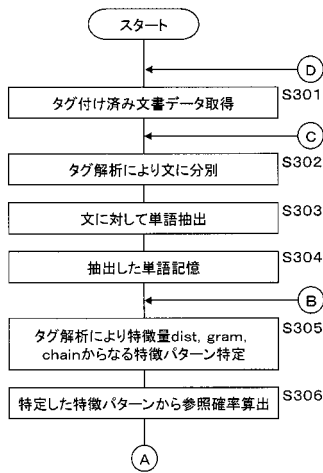
【図4】

祭とは、神霊などを祀る儀式。祭礼、祭祀とも呼ばれる。九州地方北部では、秋に行われるものに対して(お)くんちと称する場合もある。あるいは、本来の祭から派生した、催事(催し、イベント)、フェスティバルのこと。

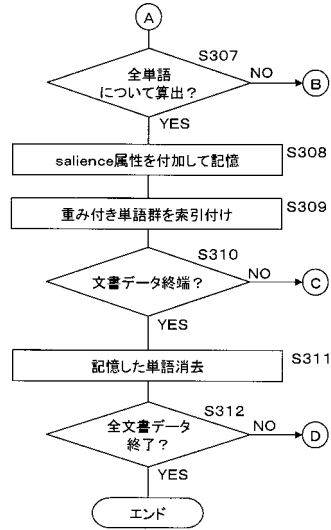
<http://ja.wikipedia.org/wiki/祭より抜粋>  
 ( <http://ja.wikipedia.org/wiki/%E7%A5%AD> )



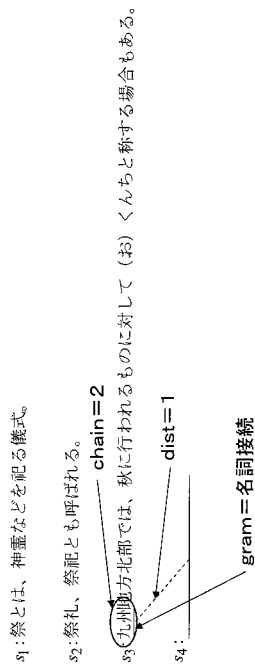
【図9】



【図10】



【図11】



【図12】

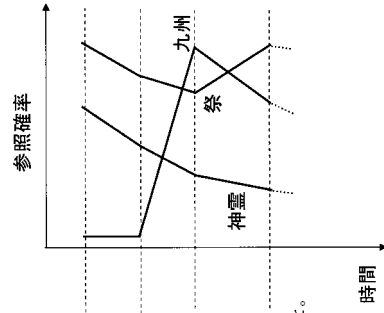
```

<su syn="f">(略)祭とは、神霊などを祀る様式。(略)</su>
<su syn="f">(略)祭礼、祭祀とも呼ばれる。(略)</su>
<su syn="f" salience="9714.0.238.23013:0.1192.9716:0.1159.13626:0.1159...">
<adp syn="f">
<n syn="f">
<placename mph="chasen;名詞+固有名詞+地域+一般::九州;キウシュウ">九州</placename>
<n mph="chasen;名詞+一般::地方;ホウフ">地方</n>
<n mph="chasen;名詞+一般::北部;ホクブ">北部</n>
</n>
<ad mph="chasen;助詞+格助詞+一般::で;デ">で</ad>
<ad mph="chasen;助詞+格助詞::は;ハ">は</ad>
</adp>
</adp>
</vp>
(略)
(略)
秋に行われるものに対して(お)くんちと称する
</vp>
<n mph="chasen;名詞+副詞可能::場合;バアイ">場合</n>
<ad mph="chasen;助詞+係助詞::も;モ">も</ad>
</adp>
<v mph="chasen;助詞::五段-う行アル+基本形:ある;アル">ある</v>
</su>
<su syn="f" salience="9714.0.238.22953:0.1836.23013:0.1192.9716:0.1159...">
(略)
あるいは、本来の祭から派生した、催事(催し、イベント)、フェスティバルのこと。
(略)
</su>
  
```

【図13】

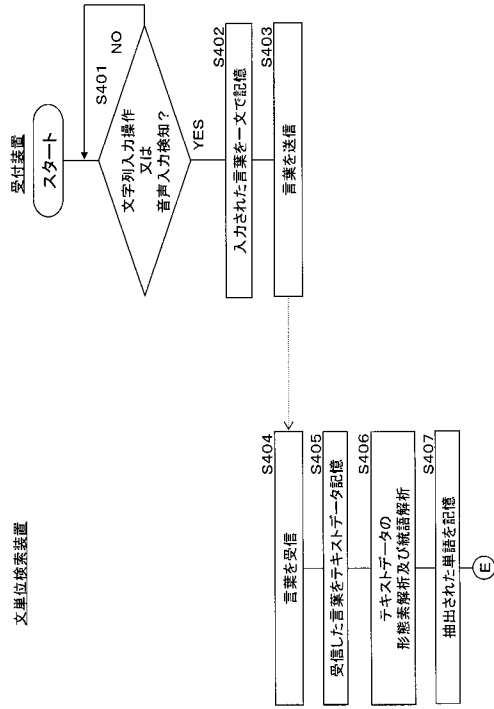
属するグループ (k-d treeノードID)	重み付き単語群	文書データ	文の位置
121	(9714:0, 238 22953:0, 1836 23013:0, 1192 9716:0, 1159 ...)	"ja.wikipedia.org/wiki/% E7%A5%AD.gda"	"/hmi/body[1]/p[2]/s[8]"
121	...	...	...
...	...	...	...
122	...	...	...
...	...	...	...

【図14】

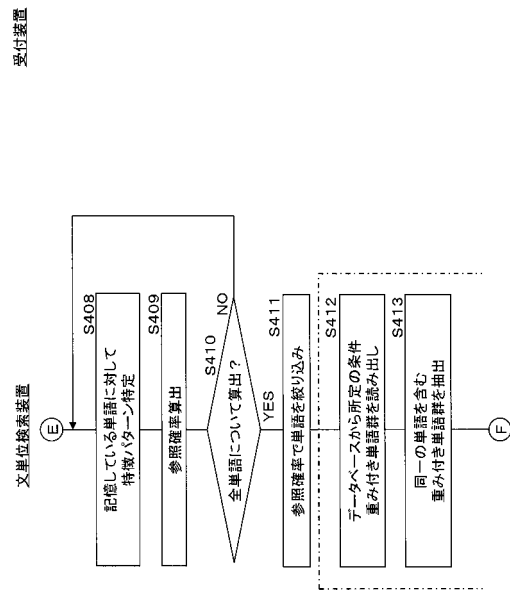


s<sub>1</sub>:祭とは、神霊などを祀る儀式。  
s<sub>2</sub>:祭礼、祭祀とも呼ばれる。  
s<sub>3</sub>:九州地方北部では、秋に行われるものに対して(お)くんちと称する場合もある。  
s<sub>4</sub>:あるいは、本来の祭から派生した、催事(催し、イベント)、フェスティバルのこと。

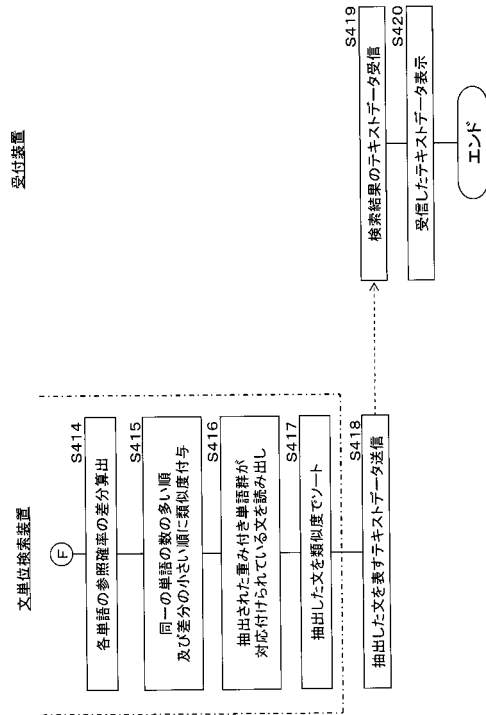
【図15】



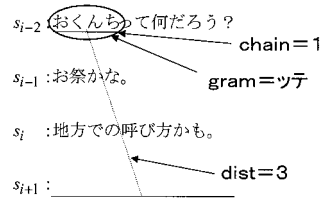
【図16】



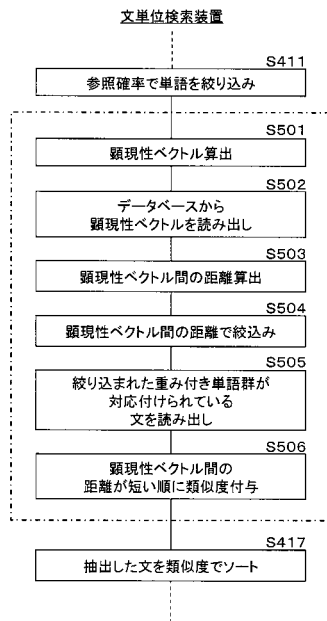
【図17】



【図18】



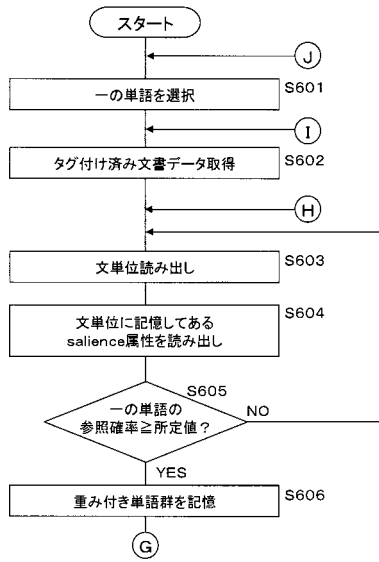
【図19】



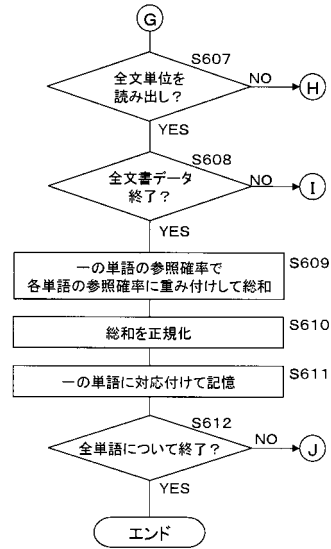
【図20】

- U<sub>1</sub>: アメリカ村も面白いですよ
- U<sub>2</sub>: それは一体どこにあるんですか
- U<sub>3</sub>: アメリカ村もミナミのちょっと横なんですけど
- U<sub>4</sub>: それは昔からある若者の町で
- U<sub>5</sub>: アメリカ村に有名な三角公園という公園があるんですね

【図 2 1】



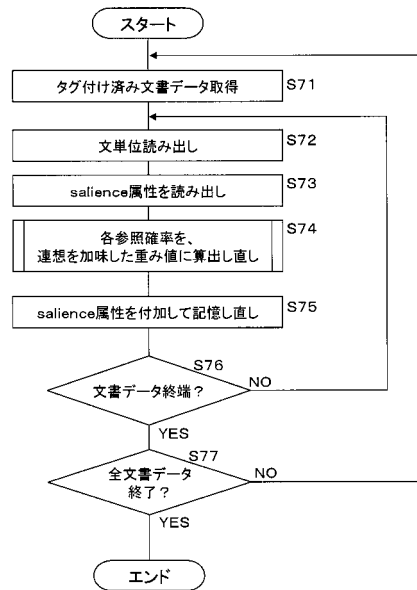
【図 2 2】



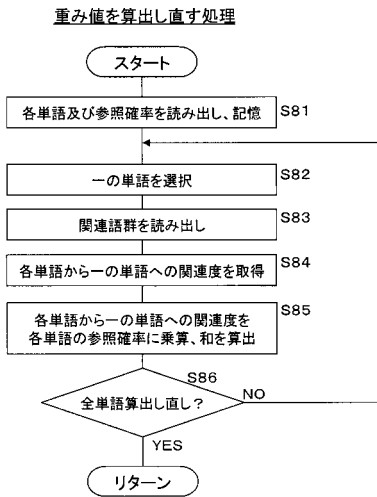
【図 2 3】

- (a)
  - GW<sub>1</sub>: (秋:0, アメリカ村:0.6, ..., 大熊座:0, 大阪:0.2, 大鹿:0, ...)
  - GW<sub>2</sub>: (秋:0.1, アメリカ村:0.3, ..., 大熊座:0, 大阪:0.4, 大鹿:0, ...)
  - GW<sub>3</sub>: (秋:0, アメリカ村:0.2, ..., 大熊座:0, 大阪:0.2, 大鹿:0, ...)
- (b)
  - GW'<sub>1</sub>: (秋:0, アメリカ村:0.36, ..., 大熊座:0, 大阪:0.12, 大鹿:0, ...)
  - GW'<sub>2</sub>: (秋:0.03, アメリカ村:0.09, ..., 大熊座:0, 大阪:0.12, 大鹿:0, ...)
  - GW'<sub>3</sub>: (秋:0, アメリカ村:0.04, ..., 大熊座:0, 大阪:0.04, 大鹿:0, ...)
- (c)
  - GW''<sub>1</sub>: (秋:0.03, アメリカ村:0.49, ..., 大熊座:0, 大阪:0.28, 大鹿:0, ...)

【図 2 4】



【図 25】



【図 26】

(a)

s<sub>1</sub>: 若者の町で。

'大阪':0.3338 '東京':0.3338 '関西人':0.2978 'ミナミ':0.201 'キタ':0.1837  
 '火':0.1637 '関西':0.1291 'アメリカ':0.1265 '心斎橋':0.1265 '町':0.121 ...

s<sub>2</sub>: アメリカ村に有名な三角公園という公園があるんですね。

'大阪':0.3208 '東京':0.3208 '関西人':0.2862 'アメリカ村':0.2021  
 'ミナミ':0.1932 '公園':0.1772 'キタ':0.1765 '火':0.1573 'アメリカ':0.1308 ...

(b)

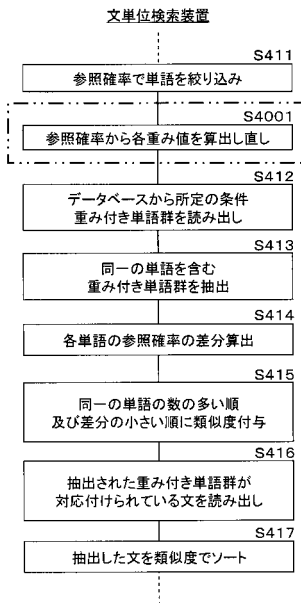
s<sub>1</sub>: 若者の町で。

'大阪':0.6229 '人':0.5345 '東京':0.3287 '自分':0.2545 '関西':0.1914  
 'アメリカ':0.1858 '昔':0.1382 '今':0.1314 '関西人':0.1261 'キタ':0.1166 ...

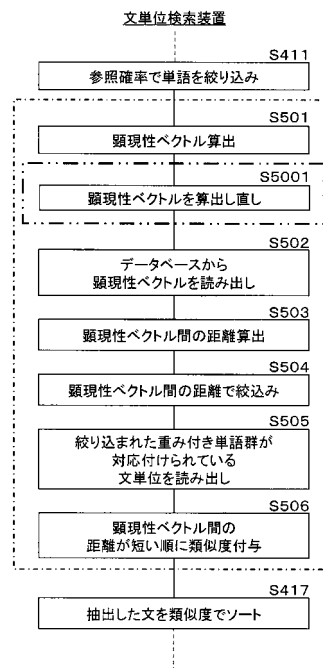
s<sub>2</sub>: アメリカ村に有名な三角公園という公園があるんですね。

'大阪':0.6675 '人':0.5625 '東京':0.3182 'アメリカ':0.2604 '自分':0.2472  
 '関西':0.1889 '公園':0.1571 '今':0.1413 '昔':0.1374 'アメリカ村':0.1339

【図 27】

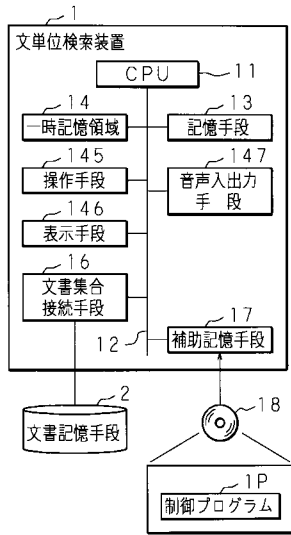


【図 28】





【図 29】



---

フロントページの続き

(72)発明者 奥乃 博

京都府京都市左京区吉田本町 京都大学大学院情報学研究科内

審査官 吉田 誠

(56)参考文献 特開平06-162092(JP,A)

特開2004-234175(JP,A)

特開2005-250762(JP,A)

西脇 正通, 関連記事を利用したテキストセグメンテーション, 情報処理学会研究報告, 日本, 社団法人情報処理学会, 2002年11月13日, Vol.2002 No.104, 79-84ページ

徳永 健伸, 言語と計算5 情報検索と言語処理, 財団法人東京大学出版会, 1999年11月25日, 第1版, 26-50ページ

(58)調査した分野(Int.Cl., DB名)

G06F 17/30