

(19) 日本国特許庁(JP)

(12) 特許公報(B2)

(11) 特許番号

特許第4849596号
(P4849596)

(45) 発行日 平成24年1月11日(2012.1.11)

(24) 登録日 平成23年10月28日(2011.10.28)

(51) Int.Cl. F I
G06F 17/30 (2006.01)
 G06F 17/30 180A
 G06F 17/30 320D
 G06F 17/30 210A

請求項の数 10 (全 62 頁)

(21) 出願番号	特願2005-354207 (P2005-354207)	(73) 特許権者	301022471
(22) 出願日	平成17年12月8日(2005.12.8)		独立行政法人情報通信研究機構
(65) 公開番号	特開2007-157006 (P2007-157006A)		東京都小金井市貫井北町4-2-1
(43) 公開日	平成19年6月21日(2007.6.21)	(74) 代理人	100094662
審査請求日	平成20年11月14日(2008.11.14)		弁理士 穂坂 和雄
		(74) 代理人	100096530
			弁理士 今村 辰夫
		(74) 代理人	100119161
			弁理士 重久 啓子
		(72) 発明者	村田 真樹
			東京都小金井市貫井北町4-2-1 独立 行政法人情報通信研究機構内
		(72) 発明者	馬 青
			東京都小金井市貫井北町4-2-1 独立 行政法人情報通信研究機構内

最終頁に続く

(54) 【発明の名称】 質問応答装置、質問応答方法および質問応答プログラム

(57) 【特許請求の範囲】

【請求項1】

第1のキーワード自体が複数のキーワードであり、かつ、第2のキーワード自体が複数のキーワードであり、第1のキーワード、第2キーワードが入力される構成であり、第1のキーワードと第2のキーワードから構成される自然言語で表現された質問データに対する解答とともに、第1のキーワードを増加して得る第3のキーワードと、第2のキーワードを増加して得る第4のキーワードとから構成される自然言語で表現された質問データに対する解答を出力する質問応答装置であって、

複数のキーワードが入力キーワードとして入力されるキーワード入力手段と、

前記入力キーワードに基づいて、前記入力キーワードの数より多いキーワードを抽出して出力キーワードとして出力するキーワード増加手段と、

前記キーワード増加手段により第1のキーワードを入力キーワードとして用いて増加して得た出力キーワードである前記第3のキーワードと、前記キーワード増加手段により第2のキーワードを入力キーワードとして用いて増加して得た出力キーワードである前記第4のキーワードとによって構成される質問に対する解答の候補である解答候補を、予め記憶された解答候補の検索対象である文書データ群から抽出する解答候補抽出手段と、

前記抽出された各解答候補が質問と対応付けられた表を解答表として出力する解答表出力手段とを備え、

前記キーワード増加手段は、

前記入力キーワードをキーワード抽出用の文書データが格納されたキーワード抽出用デ

10

20

データベースで全文検索し、前記入力キーワードのうちの複数のキーワードの検索結果において前記複数のキーワードの前後に出現する文字列をパターンとして抽出するパターン抽出手段と、

前記パターン抽出手段で抽出したパターンを前記キーワード抽出用データベースで全文検索し、前記パターンに囲まれた表現を抽出し、前記抽出した表現を出力キーワードとして出力するキーワード抽出手段とを備える、

ことを特徴とする質問応答装置。

【請求項 2】

請求項 1 に記載の質問応答装置において、

前記キーワード増加手段は、

前記入力された第 1 のキーワードに基づいて、第 3 のキーワードを出力キーワードとして出力し、前記入力された第 2 のキーワードに基づいて、第 4 のキーワードを出力キーワードとして出力し、

前記解答候補抽出手段は、予め用意された問題とその問題に対する解答の組の多数のセットを用いて、どういう問題のときにどういう解答になるかを学習し、その学習結果に基づいて、前記出力された第 3 のキーワードと第 4 のキーワードとによって構成される質問に対する解答の候補である解答候補を抽出する

ことを特徴とする質問応答装置。

【請求項 3】

請求項 1 に記載の質問応答装置において、

前記キーワード増加手段は、

前記入力された第 1 のキーワードに基づいて、第 3 のキーワードを出力キーワードとして出力し、前記入力された第 2 のキーワードに基づいて、第 4 のキーワードを出力キーワードとして出力し、

前記解答候補抽出手段は、予め記憶手段中に格納された大量の文書データ群中から前記出力された第 3 のキーワードと第 4 のキーワードを含む文書データを取り出し、取り出された文書データの言語表現から、前記大量の文書データ群中に出現する頻度を用いて、前記出力された第 3 のキーワードと第 4 のキーワードとによって構成される質問に対する解答候補を抽出する

ことを特徴とする質問応答装置。

【請求項 4】

請求項 1 に記載の質問応答装置において、

前記第 2 のキーワードに対応付けられた疑問代名詞が入力される疑問代名詞入力手段と

前記疑問代名詞入力手段により入力された疑問代名詞に基づいて、前記キーワード増加手段によって出力される出力キーワードによって構成される質問に対する解答の候補の言語表現の類型である解答タイプを推定する解答タイプ推定手段とを備え、

前記キーワード増加手段は、前記入力された第 1 のキーワードに基づいて、第 3 のキーワードを出力キーワードとして出力し、前記入力された第 2 のキーワードを出力キーワードとして出力し、

前記解答候補抽出手段は、前記解答候補の検索対象である文書データ群から、前記キーワード増加手段によって出力された第 3 のキーワードと第 2 のキーワードとを含む文書データを検索し、この検索処理で抽出された文書データから、前記解答タイプ推定手段によって推定された解答タイプに適合する言語表現を、前記第 3 のキーワードと第 2 のキーワードとによって構成される質問の解答候補として抽出する

ことを特徴とする質問応答装置。

【請求項 5】

請求項 1 に記載の質問応答装置において、

予め定められた前記第 2 のキーワードに対応付けられた疑問代名詞に基づいて、前記キーワード増加手段によって出力される出力キーワードによって構成される質問に対する解

10

20

30

40

50

答の候補の言語表現の類型である解答タイプを推定する解答タイプ推定手段とを備え、
前記キーワード増加手段は、前記入力された第1のキーワードに基づいて、第3のキーワードを出力キーワードとして出力し、前記入力された第2のキーワードを出力キーワードとして出力し、

前記解答候補抽出手段は、前記解答候補の検索対象である文書データ群から、前記キーワード増加手段によって出力された第3のキーワードと第2のキーワードとを含む文書データを検索し、この検索処理で抽出された文書データから、前記解答タイプ推定手段によって推定された解答タイプに適合する言語表現を、前記第3のキーワードと第2のキーワードとによって構成される質問の解答候補として抽出することを特徴とする質問応答装置。

10

【請求項6】

請求項1に記載の質問応答装置において、
前記キーワード増加手段は、
前記入力された第1のキーワードに基づいて、第3のキーワードを出力キーワードとして出力し、前記入力された第2のキーワードに基づいて、第4のキーワードを出力キーワードとして出力し、

前記解答候補抽出手段は、前記解答候補の検索対象である文書データ群から、前記キーワード増加手段によって出力された第3のキーワードと第4のキーワードとを含む文書データを検索し、この検索処理で抽出された文書データから、予め定められた解答タイプに適合する言語表現を、前記出力された第3のキーワードと第4のキーワードとによって構成される質問の解答候補として抽出する

20

ことを特徴とする質問応答装置。

【請求項7】

請求項1に記載の質問応答装置において、
前記キーワード増加手段によって出力される出力キーワードによって構成される質問に対する解答の候補の言語表現の類型である解答タイプであって、前記キーワード入力手段に入力された第2のキーワードに対応付けられた解答タイプが入力される解答タイプ入力手段を備え、

前記キーワード増加手段は、前記入力された第1のキーワードに基づいて、第3のキーワードを出力キーワードとして出力し、前記入力された第2のキーワードに基づいて、第4のキーワードを出力キーワードとして出力し、

30

前記第2のキーワードのうち前記出力された第4のキーワードに類似するものを、前記第4のキーワードのそれぞれについて、類似キーワードとして決定する類似キーワード決定手段を備え、

前記解答候補抽出手段は、前記解答候補の検索対象である文書データ群から、前記キーワード増加手段によって出力された第3のキーワードと第4のキーワードを含む文書データを検索し、この検索処理で抽出された文書データから、前記出力された第4のキーワードが類似する類似キーワードに対応付けられて前記解答タイプ入力手段に入力された解答タイプに適合する言語表現を、前記出力された第3のキーワードと第4のキーワードとによって構成される質問の解答候補として抽出する

40

ことを特徴とする質問応答装置。

【請求項8】

請求項7に記載の質問応答装置において、
前記類似キーワード決定手段は、
予め記憶手段内に格納された大量の文書データ群中から、前記キーワード抽出手段によって出力された第4のキーワードと共起して出現する語である共起語を抽出するとともに、前記第4のキーワードのそれぞれについて、前記抽出された各共起語と共起して前記文書データ群中に出現する回数を要素とするベクトルである共起ベクトルを求め、

各第4のキーワードについての共起ベクトルと前記キーワード入力手段に入力された第2のキーワードと同一の第4のキーワードについての共起ベクトルとの類似の度合いを求

50

め、求められた類似の度合いに基づいて決まる、前記各第4のキーワードと類似する第2のキーワードと同一の第4のキーワードを、前記類似キーワードとする

ことを特徴とする質問応答装置。

【請求項9】

第1のキーワード自体が複数のキーワードであり、かつ、第2のキーワード自体が複数のキーワードであり、第1のキーワード、第2キーワードが入力される構成であり、第1のキーワードと第2のキーワードから構成される自然言語で表現された質問データに対する解答とともに、第1のキーワードを増加して得る第3のキーワードと、第2のキーワードを増加して得る第4のキーワードとから構成される自然言語で表現された質問データに対する解答を出力する質問応答方法であって、

10

複数のキーワードで構成される第1のキーワードと複数のキーワードで構成される第2のキーワードとを入力するステップと、

入力された前記第1のキーワードと第2のキーワードに基づいて、前記第1のキーワードと第2のキーワードを構成するそれぞれのキーワードの数を増加した第3のキーワードと第4のキーワードを抽出して出力するステップと、

前記第3のキーワードと第4のキーワードを含む増加した複数の出力キーワードによって構成される質問に対する解答の候補である解答候補を、予め記憶された解答候補の検索対象である文書データ群から抽出するステップと、

前記抽出された各解答候補と質問とが対応付けられた表を解答表として出力するステップとを有し、

20

前記入力キーワードの数より多いキーワードを抽出するステップは、

前記入力キーワードをキーワード抽出用の文書データが格納されたキーワード抽出用データベースで全文検索し、前記入力キーワードのうちの複数のキーワードの検索結果において前記複数のキーワードの前後に出現する文字列をパターンとして抽出するステップと

前記抽出したパターンを前記キーワード抽出用データベースで全文検索し、前記パターンに囲まれた表現を抽出し、前記抽出した表現を出力キーワードとして出力するステップとを備える、

ことを特徴とする質問応答方法。

【請求項10】

30

第1のキーワード自体が複数のキーワードであり、かつ、第2のキーワード自体が複数のキーワードであり、第1のキーワード、第2キーワードが入力される構成であり、第1のキーワードと第2のキーワードから構成される自然言語で表現された質問データに対する解答とともに、第1のキーワードを増加して得る第3のキーワードと、第2のキーワードを増加して得る第4のキーワードとから構成される自然言語で表現された質問データに対する解答を出力する質問応答装置が備えるコンピュータに実行させるための質問応答プログラムであって、

前記コンピュータに、

複数のキーワードで構成される第1のキーワードと複数のキーワードで構成される第2のキーワードとを入力する処理と、

40

入力された前記第1のキーワードと第2のキーワードに基づいて、前記第1のキーワードと第2のキーワードを構成するそれぞれのキーワードの数を増加した第3のキーワードと第4のキーワードを抽出して出力する処理と、

前記第3のキーワードと第4のキーワードを含む増加した複数の出力キーワードによって構成される質問に対する解答の候補である解答候補を、予め記憶された解答候補の検索対象である文書データ群から抽出する処理と、

前記抽出された各解答候補と質問とが対応付けられた表を解答表として出力する処理とを実行させると共に、

前記入力キーワードの数より多いキーワードを抽出する処理は、

前記入力キーワードをキーワード抽出用の文書データが格納されたキーワード抽出用デ

50

ータベースで全文検索し、前記入力キーワードのうちの複数のキーワードの検索結果において前記複数のキーワードの前後に出現する文字列をパターンとして抽出する処理と、

前記抽出したパターンを前記キーワード抽出用データベースで全文検索し、前記パターンに囲まれた表現を抽出し、前記抽出した表現を出力キーワードとして出力する処理とを実行させる

ことを特徴とする質問応答プログラム。

【発明の詳細な説明】

【技術分野】

【0001】

本発明は、コンピュータによる自然言語処理システムとして、自然言語で表現された質問に対する解答を出力する質問応答技術に関し、特に、入力されたキーワードをキーワード抽出技術によって増加させ、増加したキーワードによって構成される複数の質問に対する解を自動的に求めて出力する質問応答装置、質問応答方法および質問応答プログラムに関する。

10

【0002】

質問応答装置とは、自然言語による質問を入力すると、その解答そのものを出力する装置である。例えば、「パーキンソン病の兆候は脳のどの部分にある細胞の死が関係していますか。」という質問を入力すると、Web、新聞記事、事典などのデータを含む大量の電子化テキストから「パーキンソン病は、中脳の黒質にあるメラニン細胞が変性し、黒質細胞内で作られる神経伝達物質のドーパミンがなくなり発病する、とされている。」といった文を探し出し、「黒質」と的確に解答を出力する。

20

【0003】

質問応答装置は、論理式やデータベースからではなく、自然言語で記述された普通の文（テキストデータ）から解答を取り出すことができるため、大量の既存の文書データを利用することができる。また、質問応答装置は、キーワードで検索された記事から使用者自らが解答を探す必要がある情報検索システムなどと異なり、解答自体を出力する。そのため、使用者は、より早く解答の情報を得ることができる。このように質問応答装置は有用であるため、より使いやすい実用的な質問応答装置の実現が期待されている。

【背景技術】

【0004】

一般的な質問応答装置（または質問応答システム）は、おおまかに、解答表現推定処理、文書検索処理、解答抽出処理という3つの処理手段で構成されている。

30

【0005】

解答表現推定処理は、入力した質問中の疑問代名詞の表現などに基づいて解答表現を推定する処理である。解答表現とは、所望される解答の言語表現の類型であって、解答となる言語表現の意味に基づいた類型（解答タイプ）、解答となる言語表現の表記に基づいた類型（解答表現タイプ）などがある。質問応答装置は、どのような質問の言語表現がどのような解答表現を要求しているかという対応関係を参照して、入力した質問の解答の解答タイプを推定する。質問応答装置は、例えば、入力した質問が「日本の面積はどのくらいですか」である場合には、所定の対応関係を参照して、質問中の「どのくらい」という表現から解答タイプは「数値表現」とであると推定する。また、質問が「日本の首相はだれですか」という場合には、質問中の「だれ」という表現から、解答タイプは「固有名詞（人名）」であると推定する。

40

【0006】

文書検索処理は、質問からキーワードを取り出し、このキーワードを用いて解答を検索する対象となっている文書データ群を検索し、解答が記述されていると考えられる文書データを抽出する処理である。質問応答装置は、例えば、入力された質問が「日本の首都はどこですか」である場合に、質問から「日本」および「首都」をキーワードとして抽出し、検索対象の文書データ群から、キーワード「日本」および「首都」を含む文書データを検索する。

50

【0007】

解答抽出処理は、文書検索処理で抽出されたキーワードを含む文書データから、推定した解答タイプに適合する言語表現を抽出し、解答として出力する処理である。質問応答装置は、例えば、文書検索処理において検索されたキーワード「日本」および「首都」を含む文書データから、解答表現推定処理において推定した解答タイプ「固有名詞（地名）」に適合する言語表現「東京」を抽出して解答とする。

【0008】

前記のような処理を行うことにより、質問応答装置は、質問「日本の首都はどこですか」に対して解答「東京」を出力する。

【0009】

なお、質問応答装置（または質問応答システム）に関する具体的な従来技術として、例えば、下記の非特許文献1に、複数の記事を使って解答の推定を行う質問応答システムにおいて、複数の記事から得られた解答の候補の得点を少しずつ減らしながら加算し、合計点が最も高い候補を解答として出力する技術について記載されている。

【非特許文献1】村田真樹，井佐原均，質問応答システムにおける遞減加点法に基づく複数記事情報の利用，情報処理学会自然言語処理研究会 2004-NL-160，2004年．九州大学．

【発明の開示】

【発明が解決しようとする課題】

【0010】

従来の質問応答装置では、検索された文書データから解答となりうる言語表現を解答候補として抽出し、抽出した解答候補それぞれの解答タイプを判定する。そして、質問から推定した解答タイプと同じか類似する解答タイプと判定した解答候補の評価を高くし、原則的には、解答タイプが同じ解答候補であって所定の評価を得たものを解答として出力する。

【0011】

しかし、従来の質問応答装置は、質問の入力によって問い合わせられた質問に対する解答のみを出力するシステムであって、問い合わせられた質問以外の質問に対する解答を出力することはできなかった。

【0012】

本発明は、上記従来技術の問題点を解決し、問い合わせられた質問に対する解答および問い合わせられた質問以外の質問に対する解答を出力する質問応答装置、質問応答方法および質問応答プログラムの提供を目的とする。

【課題を解決するための手段】

【0013】

上記課題を解決するため、本発明は、自然言語で表現された質問データに対する解答を出力する質問応答装置であって、複数のキーワードが入力キーワードとして入力されるキーワード入力手段と、前記入力キーワードに基づいて、前記入力キーワードの数より多いキーワードを抽出して出力キーワードとして出力するキーワード増加手段と、前記出力キーワードによって構成される質問に対する解答の候補である解答候補を、予め記憶された解答候補の検索対象である文書データ群から抽出する解答候補抽出手段と、前記抽出された各解答候補が質問と対応付けられた表を解答表として出力する解答表出力手段とを備えることを特徴とする。

【0014】

また、本発明は、前記の質問応答装置において、前記キーワード増加手段は、前記入力キーワードをキーワード抽出用の文書データが格納されたキーワード抽出用データベースで全文検索し、前記入力キーワードの周辺に出現したパターンを抽出するパターン抽出手段と、前記パターン抽出手段で抽出したパターンを前記キーワード抽出用データベースで全文検索し、前記パターンによって抽出される表現を抽出し、前記抽出した表現を出力キーワードとして出力するキーワード抽出手段とを備えることを特徴とする。

【0015】

また、本発明は、前記の質問応答装置において、前記キーワード増加手段は、前記入力キーワードと同じ分野の単語を、単語と単語の分野との対応情報が格納されたデータベースから抽出し、出力キーワードとして出力することを特徴とする。

【0016】

また、本発明は、前記の質問応答装置において、前記キーワード増加手段は、予めデータベース中に記憶された、意味的類似による単語の分類情報であるシソーラスデータに基づいて、前記入力された入力キーワードと、前記シソーラスデータ中の単語との類似度を算出する類似度算出手段と、前記算出された類似度の大きさに基づいてキーワードを抽出し、出力キーワードとして出力するキーワード抽出手段とを備えることを特徴とする。

【0017】

また、本発明は、前記の質問応答装置において、前記キーワード入力手段には、前記入力キーワードとして第1のキーワードと第2のキーワードとが入力され、前記キーワード増加手段は、前記入力された第1のキーワードに基づいて、第3のキーワードを出力キーワードとして出力し、前記入力された第2のキーワードに基づいて、第4のキーワードを出力キーワードとして出力し、前記解答候補抽出手段は、予め用意された問題とその問題に対する解答の組の多数のセットを用いて、どういう問題のときにどういう解答になるかを学習し、その学習結果に基づいて、前記出力された第3のキーワードと第4のキーワードとによって構成される質問に対する解答の候補である解答候補を抽出することを特徴とする。

【0018】

また、本発明は、前記の質問応答装置において、前記キーワード入力手段には、前記入力キーワードとして第1のキーワードと第2のキーワードとが入力され、前記キーワード増加手段は、前記入力された第1のキーワードに基づいて、第3のキーワードを出力キーワードとして出力し、前記入力された第2のキーワードに基づいて、第4のキーワードを出力キーワードとして出力し、前記解答候補抽出手段は、予め記憶手段中に格納された大量の文書データ群中から前記出力された第3のキーワードと第4のキーワードを含む文書データを取り出し、取り出された文書データの言語表現から、前記大量の文書データ群中出现する頻度を用いて、前記出力された第3のキーワードと第4のキーワードとによって構成される質問に対する解答候補を抽出することを特徴とする。

【0019】

また、本発明は、前記の質問応答装置において、前記キーワード入力手段には、前記入力キーワードとして第1のキーワードと第2のキーワードとが入力され、前記第2のキーワードに対応付けられた疑問代名詞が入力される疑問代名詞入力手段と、前記疑問代名詞入力手段に入力された疑問代名詞に基づいて、前記キーワード増加手段によって出力される出力キーワードによって構成される質問に対する解答の候補の言語表現の類型である解答タイプを推定する解答タイプ推定手段とを備え、前記キーワード増加手段は、前記入力された第1のキーワードに基づいて、第3のキーワードを出力キーワードとして出力し、前記入力された第2のキーワードを出力キーワードとして出力し、前記解答候補抽出手段は、前記解答候補の検索対象である文書データ群から、前記キーワード増加手段によって出力された第3のキーワードと第2のキーワードとを含む文書データを検索し、この検索処理で抽出された文書データから、前記解答タイプ推定手段によって推定された解答タイプに適合する言語表現を、前記第3のキーワードと第2のキーワードとによって構成される質問の解答候補として抽出することを特徴とする。

【0020】

また、本発明は、前記の質問応答装置において、前記キーワード入力手段には、前記入力キーワードとして第1のキーワードと第2のキーワードとが入力され、予め定められた前記第2のキーワードに対応付けられた疑問代名詞に基づいて、前記キーワード増加手段によって出力される出力キーワードによって構成される質問に対する解答の候補の言語表現の類型である解答タイプを推定する解答タイプ推定手段とを備え、前記キーワード増加手段は、前記入力された第1のキーワードに基づいて、第3のキーワードを出力キーワー

10

20

30

40

50

ドとして出力し、前記入力された第2のキーワードを出力キーワードとして出力し、前記解答候補抽出手段は、前記解答候補の検索対象である文書データ群から、前記キーワード増加手段によって出力された第3のキーワードと第2のキーワードとを含む文書データを検索し、この検索処理で抽出された文書データから、前記解答タイプ推定手段によって推定された解答タイプに適合する言語表現を、前記第3のキーワードと第2のキーワードとによって構成される質問の解答候補として抽出することを特徴とする。

【0021】

また、本発明は、前記の質問応答装置において、前記キーワード入力手段には、前記入力キーワードとして第1のキーワードと第2のキーワードとが入力され、疑問代名詞が入力される疑問代名詞入力手段と、前記疑問代名詞入力手段に入力された疑問代名詞に基づいて、前記キーワード増加手段によって出力される出力キーワードによって構成される質問に対する解答の候補の言語表現の類型である解答タイプを推定する解答タイプ推定手段とを備え、前記キーワード増加手段は、前記入力された第1のキーワードに基づいて、第3のキーワードを出力キーワードとして出力し、前記入力された第2のキーワードに基づいて、第4のキーワードを出力キーワードとして出力し、前記解答候補抽出手段は、前記解答候補の検索対象である文書データ群から、前記キーワード増加手段によって出力された第3のキーワードと第2のキーワードとを含む文書データを検索し、この検索処理で抽出された文書データから、前記解答タイプ推定手段によって推定された解答タイプに適合する言語表現を、前記第3のキーワードと第4のキーワードとによって構成される質問の解答候補として抽出することを特徴とする。

【0022】

また、本発明は、前記の質問応答装置において、前記キーワード入力手段には、前記入力キーワードとして第1のキーワードと第2のキーワードとが入力され、予め定められた疑問代名詞に基づいて、前記キーワード増加手段によって出力される出力キーワードによって構成される質問に対する解答の候補の言語表現の類型である解答タイプを推定する解答タイプ推定手段とを備え、前記キーワード増加手段は、前記入力された第1のキーワードに基づいて、第3のキーワードを出力キーワードとして出力し、前記入力された第2のキーワードに基づいて、第4のキーワードを出力キーワードとして出力し、前記解答候補抽出手段は、前記解答候補の検索対象である文書データ群から、前記キーワード増加手段によって出力された第3のキーワードと第2のキーワードとを含む文書データを検索し、この検索処理で抽出された文書データから、前記解答タイプ推定手段によって推定された解答タイプに適合する言語表現を、前記第3のキーワードと第4のキーワードとによって構成される質問の解答候補として抽出することを特徴とする。

【0023】

また、本発明は、前記の質問応答装置において、前記キーワード増加手段によって出力される出力キーワードによって構成される質問に対する解答の候補の言語表現の類型である解答タイプが入力される解答タイプ入力手段を備え、前記キーワード入力手段には、前記入力キーワードとして第1のキーワードと第2のキーワードとが入力され、前記キーワード増加手段は、前記入力された第1のキーワードに基づいて、第3のキーワードを出力キーワードとして出力し、前記入力された第2のキーワードに基づいて、第4のキーワードを出力キーワードとして出力し、前記解答候補抽出手段は、前記解答候補の検索対象である文書データ群から、前記キーワード増加手段によって出力された第3のキーワードと第4のキーワードとを含む文書データを検索し、この検索処理で抽出された文書データから、前記解答タイプ入力手段に入力された解答タイプに適合する言語表現を、前記出力された第3のキーワードと第4のキーワードとによって構成される質問の解答候補として抽出することを特徴とする。

【0024】

また、本発明は、前記の質問応答装置において、前記キーワード入力手段には、前記入力キーワードとして第1のキーワードと第2のキーワードとが入力され、前記キーワード増加手段は、前記入力された第1のキーワードに基づいて、第3のキーワードを出力キー

10

20

30

40

50

ワードとして出力し、前記入力された第2のキーワードに基づいて、第4のキーワードを出力キーワードとして出力し、前記解答候補抽出手段は、前記解答候補の検索対象である文書データ群から、前記キーワード増加手段によって出力された第3のキーワードと第4のキーワードとを含む文書データを検索し、この検索処理で抽出された文書データから、予め定められた解答タイプに適合する言語表現を、前記出力された第3のキーワードと第4のキーワードとによって構成される質問の解答候補として抽出することを特徴とする。

【0025】

また、本発明は、前記の質問応答装置において、前記キーワード入力手段には、前記入力キーワードとして第1のキーワードと第2のキーワードとが入力され、前記キーワード増加手段によって出力される出力キーワードによって構成される質問に対する解答の候補の言語表現の類型である解答タイプであって、前記キーワード入力手段に入力された第2のキーワードに対応付けられた解答タイプが入力される解答タイプ入力手段を備え、前記キーワード増加手段は、前記入力された第1のキーワードに基づいて、第3のキーワードを出力キーワードとして出力し、前記入力された第2のキーワードを出力キーワードとして出力し、前記解答候補抽出手段は、前記解答候補の検索対象である文書データ群から、前記キーワード増加手段によって出力された第3のキーワードと第2のキーワードとを含む文書データを検索し、この検索処理で抽出された文書データから、前記解答タイプ入力手段に入力された解答タイプに適合する言語表現を、前記出力された第3のキーワードと第2のキーワードとによって構成される質問の解答候補として抽出することを特徴とする。

【0026】

また、本発明は、前記の質問応答装置において、前記キーワード入力手段には、前記入力キーワードとして第1のキーワードと第2のキーワードとが入力され、前記キーワード増加手段は、前記入力された第1のキーワードに基づいて、第3のキーワードを出力キーワードとして出力し、前記入力された第2のキーワードを出力キーワードとして出力し、前記解答候補抽出手段は、前記解答候補の検索対象である文書データ群から、前記キーワード増加手段によって出力された第3のキーワードと第2のキーワードとを含む文書データを検索し、この検索処理で抽出された文書データから、予め定められた、前記第2のキーワードに対応付けられた解答タイプに適合する言語表現を、前記出力された第3のキーワードと第2のキーワードとによって構成される質問の解答候補として抽出することを特徴とする。

【0027】

また、本発明は、前記の質問応答装置において、前記キーワード入力手段には、前記入力キーワードとして第1のキーワードと第2のキーワードとが入力され、前記キーワード増加手段によって出力される出力キーワードによって構成される質問に対する解答の候補の言語表現の類型である解答タイプであって、前記キーワード入力手段に入力された第2のキーワードに対応付けられた解答タイプが入力される解答タイプ入力手段を備え、前記キーワード増加手段は、前記入力された第1のキーワードに基づいて、第3のキーワードを出力キーワードとして出力し、前記入力された第2のキーワードに基づいて、第4のキーワードを出力キーワードとして出力し、前記第2のキーワードのうち前記出力された第4のキーワードに類似するものを、前記第4のキーワードのそれぞれについて、類似キーワードとして決定する類似キーワード決定手段を備え、前記解答候補抽出手段は、前記解答候補の検索対象である文書データ群から、前記キーワード増加手段によって出力された第3のキーワードと第4のキーワードを含む文書データを検索し、この検索処理で抽出された文書データから、前記出力された第4のキーワードが類似する類似キーワードに対応付けられて前記解答タイプ入力手段に入力された解答タイプに適合する言語表現を、前記出力された第3のキーワードと第4のキーワードとによって構成される質問の解答候補として抽出することを特徴とする。

【0028】

また、本発明は、前記の質問応答装置において、前記キーワード入力手段には、前記入

10

20

30

40

50

力キーワードとして第1のキーワードと第2のキーワードとが入力され、前記キーワード増加手段は、前記入力された第1のキーワードに基づいて、第3のキーワードを出力キーワードとして出力し、前記入力された第2のキーワードに基づいて、第4のキーワードを出力キーワードとして出力し、前記第2のキーワードのうち前記出力された第4のキーワードに類似するものを、前記第4のキーワードのそれぞれについて、類似キーワードとして決定する類似キーワード決定手段を備え、前記解答候補抽出手段は、前記解答候補の検索対象である文書データ群から、前記キーワード増加手段によって出力された第3のキーワードと第4のキーワードを含む文書データを検索し、この検索処理で抽出された文書データから、前記出力された第4のキーワードが類似する類似キーワードに予め対応付けられた解答タイプに適合する言語表現を、前記出力された第3のキーワードと第4のキーワードとによって構成される質問の解答候補として抽出することを特徴とする。

10

【0029】

また、本発明は、前記の質問応答装置において、前記キーワード入力手段には、前記入力キーワードとして第1のキーワードと第2のキーワードとが入力され、前記第2のキーワードに対応付けられた疑問代名詞が入力される疑問代名詞入力手段と、前記疑問代名詞入力手段に入力された疑問代名詞に基づいて、前記キーワード増加手段によって出力される出力キーワードによって構成される質問に対する解答の候補の言語表現のタイプである解答タイプを推定する解答タイプ推定手段とを備え、前記キーワード増加手段は、前記入力された第1のキーワードに基づいて、第3のキーワードを出力キーワードとして出力し、前記入力された第2のキーワードに基づいて、第4のキーワードを出力キーワードとして出力し、前記第2のキーワードのうち前記出力された第4のキーワードに類似するものを、前記第4のキーワードのそれぞれについて、類似キーワードとして決定する類似キーワード決定手段を備え、前記解答候補抽出手段は、前記解答候補の検索対象である文書データ群から、前記キーワード増加手段によって出力された第3のキーワードと第4のキーワードを含む文書データを検索し、この検索処理で抽出された文書データから、前記出力された第4のキーワードが類似する類似キーワードに対応付けられて前記疑問代名詞入力手段に入力された疑問代名詞に基づいて解答タイプ推定手段が推定した解答タイプに適合する言語表現を、前記出力された第3のキーワードと第4のキーワードとによって構成される質問の解答候補として抽出することを特徴とする。

20

【0030】

また、本発明は、前記の質問応答装置において、前記キーワード入力手段には、前記入力キーワードとして第1のキーワードと第2のキーワードとが入力され、予め定められた、前記第2のキーワードに対応付けられた疑問代名詞に基づいて、前記キーワード増加手段によって出力される出力キーワードによって構成される質問に対する解答の候補の言語表現のタイプである解答タイプを推定する解答タイプ推定手段とを備え、前記キーワード増加手段は、前記入力された第1のキーワードに基づいて、第3のキーワードを出力キーワードとして出力し、前記入力された第2のキーワードに基づいて、第4のキーワードを出力キーワードとして出力し、前記第2のキーワードのうち前記出力された第4のキーワードに類似するものを、前記第4のキーワードのそれぞれについて、類似キーワードとして決定する類似キーワード決定手段を備え、前記解答候補抽出手段は、前記解答候補の検索対象である文書データ群から、前記キーワード増加手段によって出力された第3のキーワードと第4のキーワードを含む文書データを検索し、この検索処理で抽出された文書データから、前記出力された第4のキーワードが類似する類似キーワードに対応付けられた疑問代名詞に基づいて解答タイプ推定手段が推定した解答タイプに適合する言語表現を、前記出力された第3のキーワードと第4のキーワードとによって構成される質問の解答候補として抽出することを特徴とする。

30

40

【0031】

また、本発明は、前記の質問応答装置において、前記類似キーワード決定手段は、予め記憶手段内に格納された大量の文書データ群中から、前記キーワード抽出手段によって出力された第4のキーワードと共起して出現する語である共起語を抽出するとともに、前記

50

第4のキーワードのそれぞれについて、前記抽出された各共起語と共起して前記文書データ群中出现する回数を要素とするベクトルである共起ベクトルを求め、各第4のキーワードについての共起ベクトルと前記キーワード入力手段に入力された第2のキーワードと同一の第4のキーワードについての共起ベクトルとの類似の度合いを求め、求められた類似の度合いに基づいて決まる、前記各第4のキーワードと類似する第2のキーワードと同一の第4のキーワードを、前記類似キーワードとすることを特徴とする。

【0032】

また、本発明は、前記の質問応答装置において、前記類似キーワード決定手段は、予めデータベース中に記憶された、意味的類似による単語の分類情報であるシソーラスデータに基づいて、前記キーワード増加手段によって出力された第4のキーワード毎に、前記第4のキーワードと同一の単語と、前記キーワード入力手段に入力された第2のキーワードと同一の単語との類似度を算出する類似度算出手段と、前記算出された類似度の大きさに基づいて決まる、前記第4のキーワードと類似する第2のキーワードを、前記類似キーワードとすることを特徴とする。

10

【0033】

また、本発明は、自然言語で表現された質問データに対する解答を出力する質問応答方法であって、複数のキーワードを入力キーワードとして入力するステップと、前記入力キーワードに基づいて、前記入力キーワードの数より多いキーワードを抽出して出力キーワードとして出力するステップと、前記出力キーワードによって構成される質問に対する解答の候補である解答候補を、予め記憶された解答候補の検索対象である文書データ群から抽出するステップと、前記抽出された各解答候補が質問と対応付けられた表を解答表として出力するステップとを有することを特徴とする。

20

【0034】

また、本発明は、自然言語で表現された質問データに対する解答を出力する質問応答装置が備えるコンピュータに実行させるためのプログラムであって、前記コンピュータに、複数のキーワードを入力キーワードとして入力する処理と、前記入力キーワードに基づいて、前記入力キーワードの数より多いキーワードを抽出して出力キーワードとして出力する処理と、前記出力キーワードによって構成される質問に対する解答の候補である解答候補を、予め記憶された解答候補の検索対象である文書データ群から抽出する処理と、前記抽出された各解答候補が質問と対応付けられた表を解答表として出力する処理とを実行させるための質問応答プログラムである。

30

【発明の効果】

【0035】

本発明の質問応答装置によれば、問い合わせられた質問に対する解答だけでなく、問い合わせられた質問以外の質問に対する解答を、各質問に対応付けた形式で出力することが可能となる。すなわち、本発明の質問応答装置によれば、ユーザは、解答を知りたいジャンルのキーワードを少数入力するだけで、入力されたキーワードに基づいて増加したキーワードによって構成される多数の質問に対する解答を自動的に得ることができる。

【0036】

例えば、本発明の質問応答装置によれば、ユーザが第1のキーワードと第2のキーワードとを入力すると、第1のキーワードに基づいて、第1のキーワードの数より多い第3のキーワードが抽出されるとともに、第2のキーワードに基づいて、第2のキーワードの数より多い第4のキーワードが抽出され、抽出された第3のキーワードと第4のキーワードに基づいて構成される質問に対する解答を機械学習の手法を用いて自動的に出力することが可能となる。

40

【0037】

また、例えば、本発明の質問応答装置によれば、ユーザが第1のキーワードと第2のキーワードと、第2のキーワードに対応付けられた疑問代名詞とを入力すると、第1のキーワードの数より多い第3のキーワードが抽出されるとともに、上記入力された疑問代名詞に基づいて解答タイプが推定され、第3のキーワードと第2のキーワードと疑問代名詞に

50

基づいて構成される質問に対する解答を、上記推定された解答タイプを用いて自動的に出力することが可能となる。

【0038】

また、例えば、本発明の質問応答装置によれば、ユーザが第1のキーワードと第2のキーワードと、解答タイプとを入力すると、第1のキーワードに基づいて、第1のキーワードの数より多い第3のキーワードが抽出されるとともに、第2のキーワードに基づいて、第2のキーワードの数より多い第4のキーワードが抽出され、抽出された第3のキーワードと第4のキーワードに基づいて構成される質問に対する解答を、上記入力された解答タイプを用いて自動的に出力することが可能となる。

【0039】

また、例えば、本発明の質問応答装置によれば、ユーザが第1のキーワードと第2のキーワードと、第2のキーワードに対応付けられた解答タイプとを入力すると、第1のキーワードに基づいて、第1のキーワードの数より多い第3のキーワードが抽出されるとともに、第2のキーワードに基づいて、第2のキーワードの数より多い第4のキーワードが抽出され、さらに、抽出された第4のキーワードに類似する第2のキーワード(と同一の第4のキーワード)が類似キーワードとして決定される。そして、抽出された第3のキーワードと第4のキーワードに基づいて構成される質問に対する解答を、上記決定された類似キーワードに対応付けられた解答タイプを用いて自動的に出力することが可能となる。

【発明を実施するための最良の形態】

【0040】

まず、本発明の実施の形態の説明の前に、上記非特許文献1に記載された技術について説明する。非特許文献1では、質問応答システムにおける逓減加点法に基づく複数記事情報の利用について記載されている。以下に非特許文献1の記載内容について説明する。

【0041】

質問応答システムは、与えられた質問に対してその答えを出力するシステムのこと、例えば、「日本の首都はどこですか」という質問文が与えられると、「東京は日本の首都で、その国の最も大きく重要な都市であり、東京は日本の47都道府県のうちの一つである。」という文をウェブや新聞記事などの電子テキストから探し出し、「東京」と答える。質問応答システムは、情報検索の代りとして重要になるだろうし、また将来の人工知能システムの基本要素にもなるであろう重要なものである。

【0042】

非特許文献1では、質問応答システムの精度向上のために、複数の記事から得た解の候補の得点を減らしながら加点する新しい方法を提案している。この方法を逓減加点法と呼ぶ。

【0043】

質問の答えが複数の記事で見つかることは多く、そのような場合は、複数の記事を使って答えを推定した方が一つの記事を使って推定するよりも良い答えを得ることができると思われるので、複数の記事から得た解の候補の得点を加算することで、複数の記事の情報を活用する手法が考えられる。しかし、ただ単純に得点を加算するだけではシステムの性能を下げる場合がある。

【0044】

そこで、非特許文献1では、この単純に加算する際に生じる問題に対処するために、得点の加算の際に得点を減らしながら加算する手法を用いる。より具体的に言うと、非特許文献1の方法では、 i 番目の解の候補の得点には $k^{(i-1)}$ の重みをかけておいて、その後で得点を加算する。最終的な答えは合計得点により判断する。例えば、「東京」が三つの記事から解の候補として抽出され、それらの得点が26、21、20であり、 k が0.3であったとする。この場合、「東京」の合計得点は、 34.1 となる(= $26 + 21 \times 0.3 + 20 \times 0.3^2$)。このような方法でそれぞれの候補の得点を計算し、最も高い合計得点を持つ候補を解とする。

【0045】

10

20

30

40

50

次に、非特許文献1における複数記事の利用における逓減加点法の利用について詳細に説明する。「日本の首都はどこですか」という質問文が与えられたとする。このとき、得るべき答えは「東京」である。一般的な質問応答システムは、図21のように、解の候補と得点をリストとして出力でき、また、解の候補を取り出した記事を指し示す記事番号も出力することができる。なお、図中に示す順位は、得点の大きさの順位を示す。

【0046】

図21に示すリストの例だと、最も得点の大きい候補は「京都」であり、誤った解を出力することになる。解の候補の得点を単純に加算する方法は、すでに提案されている。図21に示すリストを用いると、解の候補の得点を単純に加算する方法によれば、図22に示す結果を得る。

【0047】

図22では、「東京」の得点が一番順位が高く、システムは、正しく「東京」を解として出力することができる。この、解の候補の得点を単純に加算する方法は、複数の記事の情報を利用することで正しい解を得ることができた。しかし、この方法には、高頻度の解の候補を取り出しやすいという問題がある。これは、特に性能が高いシステムで深刻な問題である。もともと性能が高いシステムでは、システムの出力した元の得点の方が単純に加算した得点よりも信頼できる場合が多く、単純に加算する方法は、しばしばシステムの性能を劣化させることになる。

【0048】

この問題に対処するために、非特許文献1の技術は、得点を減らしながら加算する新しい方法を提案している。解の候補の得点を単純に加算する代わりに、得点を減らす重みをつけて得点を加算するのである。この方法は、高頻度語を取り出し易いという悪い効果を減じ、なおかつシステムの性能を向上させる効果を持つ。

【0049】

この、非特許文献1で提案する方法の有効性を示す例をあげる。「日本の首都は西暦1000年の時はどこでしたか。」と質問が与えられ、システムは図23に示す結果を出力したとする。図23に示すように、「京都」の得点が一番高い。ここで、上記質問に対する正解は「京都」であり、解の候補の得点を単純に加算しなければ、このシステムは正解を出力している。しかし、単純に加算する方法を用いると、その結果は図24に示す表のようになり、間違った解の「東京」をシステムの解としてしまう。

【0050】

ここで、得点を減らしながら加算する非特許文献1の新しい方法を利用してみる。ここでは、細かいシステムの仕様として、 i 番目の候補の得点に $0.3^{(i-1)}$ を乗じることとする。その場合、「東京」の得点は 2.8 であり ($= 2.1 + 1.8 \times 0.3 + 1.5 \times 0.3^2 + 1.4 \times 0.3^3$)、システムの出力結果は、図25に示す表のようになり、「京都」の得点が一番高いので、正解の「京都」を解として正しく出力することができる。すなわち、非特許文献1で提案する方法は、最初の例(「日本の首都はどこですか」という質問文が与えられた場合)でも正しい解を得ることができる。最初の例に適用すると、「東京」の得点は 4.3 となり ($= 3.2 + 2.8 \times 0.3 + 2.5 \times 0.3^2 + 2.4 \times 0.3^3$)、出力結果は図26に示す表のようになり、「東京」が最も高い得点となり、解として正しく出力される。

【0051】

得点を減らしながら加算する非特許文献1に記載された方法は、高頻度の解の候補を取り出しやすい欠点を減じながら、なおかつ複数記事の情報を利用し精度向上を実現できるものである。

【0052】

非特許文献1に記載された質問応答システムは、以下の三つの基本要素からなる。

1. 解表現の推定

質問応答システムは、疑問代名詞の表現などに基づいて解表現(解がどのような言語表現か)を推定する。例えば、入力の問題文が「日本の面積はどのくらいですか」だとする

10

20

30

40

50

と、「どのくらい」という表現から、解表現は数値表現であろうと推測する。

2. 文書検索

質問応答システムは、質問文からキーワードを取り出し、これらのキーワードを用いて文書を検索する。この検索により、解が書いてありそうな文書群を集めることになる。例えば、入力された質問文が、「日本の面積はどのくらいですか」とすると、「日本」、「面積」がキーワードとして抽出され、これらを含む文書を検索することになる。

3. 解の抽出

質問応答システムは、解が書いてありそうな文書群から、推定した解表現に適合する言語表現を抽出し、それを解として出力する。例えば、入力された質問文が、「日本の面積はどのくらいですか」とすると、文書検索で検索した「日本」、「面積」を含む文書群から、解表現として推定した数値表現にあたる言語表現を解として抽出する。

【0053】

以下に、非特許文献1で提案する技術について、詳細に説明する。

【0054】

(解表現の推定)

人手で作成したヒューリスティックルールを使って解表現を推定する。16個のルールを作成する。そのいくつかを以下に示す。

- ・質問文に「誰」という表現がある場合、解表現は人名である。
- ・質問文に「いつ」という表現がある場合、解表現は時間表現である。
- ・質問文に「どのくらいの」という表現がある場合、解表現は数値表現である。

【0055】

(文書検索)

文書検索のためのキーワードは、公知のキーワード抽出ツールであるChasenにより取り出し、付属語などはキーワードから除外する。文書検索は以下のように行なう。

【0056】

まず、以下の式で文書検索を行ない、上位 k_{dr1} 個の記事を取り出す。

【0057】

【数1】

$$Score(d) = \sum_{\text{term } t} \left(\frac{tf(d,t)}{tf(d,t) + k_t} \frac{length(d) + k_+}{\Delta + k_+} \times \log \frac{N}{df(t)} \right) \text{式(1)}$$

【0058】

ただし、 d は記事で、 t は質問文から取り出したキーワードで、 $tf(d, t)$ は、記事 d に出現するキーワード t の頻度で、 $df(t)$ はキーワード t が出現する頻度で、 N は記事の総数で、 $length(d)$ は記事 d の長さで、 Δ は記事長の平均である。 k_t と k_+ は実験で定める定数である。この式は、ロバートソンのOkapiウェイトイング(例えば、下記の文献(1)、文献(2)参照)の式に基づくもので、情報検索でよく用いられる式である(例えば、下記の文献(3)、文献(4)参照)。但し、質問応答では多くの種類のキーワードがマッチすることが重要なので、 k_t の値としては大きな値を用いる。

【0059】

文献(1): S.E. Robertson and S.Walker, Some simple effective approximations to the 2-poisson model for probabilistic weighted retrieval, Proceedings of the Seventeenth Annual International ACM SIGIR Conference on Research and Development in Information Retrieval, (1994).

文献(2): S.E. Robertson, S.Walker, S.Jones, M.M. Hancock-Beaulieu, and M.Gatford, Okapi at trec-3, TREC-3, (1994).

文献(3): 村田真樹, 内元清貴, 小作浩美, 馬青, 内山将夫, 井佐原均, 位置情報と

10

20

30

40

50

分野情報を用いた情報検索，言語処理学会誌，Vol. 7，No. 2 (2000)。

文献(4)：Masaki Murata, Qing Ma, and Hitoshi Isahara, High performance information retrieval using many characteristics and many techniques, Proceedings of the Third NTCIR Workshop (CLIR), (2002)。

次に、以下の式で記事をリランキングし、上位 k_{dr2} 個の記事を取り出す。

【0060】

【数2】

$$Score(d) = \max_{t1 \in T} \sum_{t2 \in T3} w_{dr2}(t2) \log \frac{N}{2dist(t1, t2) * df(t2)} \quad 10$$

式(2)

$$T3 = \{t | t \in T, 2dist(t1, t) \frac{df(t)}{N} \leq 1\} \quad 式(3)$$

【0061】

ただし、 T はキーワードの集合で、 $dist(t1, t2)$ はキーワード $t1$ と $t2$ の間の距離で、便宜上 $t1 = t2$ のとき $dist(t1, t2) = 0.5$ としている。 w_{dr2} は $t2$ の関数で実験により定められる。 20

【0062】

一般には、質問応答システムでは質問文から取り出した複数のキーワードが近くに出現することを保証するために、記事を段落などの小さい単位に分割するが、非特許文献1のシステムでは、上記の、リランキングによりキーワードが近くにある場合に得点をあげる式を用いるので、記事を分割する必要がなく、記事をそのまま文書検索に使えるのである。この文書検索では、上位20記事を取り出し、それを次の解の抽出で利用する。

【0063】

(解の抽出)

文書検索で得た記事から、名詞、未知語連続を取り出し、それらを解の候補とする。それぞれの候補には、解の候補とキーワードの近さに基づく得点 $Score_{near}(c)$ と解表現の意味制約を満足しているか否かに基づく $Score_{sem}(c)$ の二つの得点を与え、その合計点が最も大きい候補を解とする。 30

【0064】

$Score_{near}(c)$ は、以下の式で与えられる。

【0065】

【数3】

$$Score_{near}(c) = \sum_{t2 \in T3} w_{dr2}(t2) \log \frac{N}{2dist(c, t2) * df(t2)} \quad 40$$

式(4)

$$T3 = \{t | t \in T, 2dist(c, t) \frac{df(t)}{N} \leq 1\} \quad 式(5)$$

【0066】

ただし、 c は解の候補であり、 w_{dr2} は実験で定められる関数である。

【0067】

解表現の意味制約に基づく得点 $Score_{sem}(c)$ は、人手で作成した規則により与えられる。非特許文献 1 では、45 の規則を作成した。そのいくつかを以下に示す。

- ・推定した解表現（人名や地名など）と一致する候補に 1000 を与える。解の候補が人名か地名かと特定する方法には、例えば SVM に基づく固有表現抽出技術を利用する。固有表現抽出技術の例については、後述する。

- ・解表現が「国名」の場合に解の候補が国名のときに 1000 を与える。

- ・質問文が「何 + 名詞 X」の場合、名詞 X を最後に持つ候補に 1000 を与える。

【0068】

非特許文献 1 における実験では、以下の得点加算法を利用している。

(1) オリジナル法

得点の加算を行わない方法。

(2) 単純加算法

複数の記事から取り出した解の候補の得点を加算し、その得点をそのまま加算した合計得点に基づき解を出力する。

(3) 遞減加算法

複数記事から取り出した候補の得点を加算する。この方法は、 i 番目の候補の得点には $k^{(i-1)}$ の値を乗じてから得点を加算する。すなわち、加算結果は、以下の式で表される。

【0069】

【数 4】

$$Score_{decreased} = \sum_{1 \leq i \leq n} k^{i-1} score_{original}(i) \quad \text{式(6)}$$

【0070】

ただし、 $Score_{decreased}$ は、最終的な加算後の値の 1000 より下の桁の数字で、 $score_{original}(i)$ は、元の値の 1000 より下の桁の数字である。 n は 1000 より上の桁で同じ数字を持つ複数の記事から得られた同じ解の候補の出現回数である。 k は実験で定める定数である。

(4) 融合法

この方法は、オリジナル法、単純加算法、遞減加算法の組み合わせである。この方法はまず学習データでこれらの方法のうちの方法が最も良い精度を出すかを調べて、最も精度の高かった方法を利用して問題を解く。

【0071】

この方法には融合による精度向上という効果と、公平な評価ができるという効果がある。

【0072】

以下に、本発明の実施の形態について、図を用いて説明する。図 1 は、本発明の第 1 の実施の形態における質問応答装置の構成の一例を示す図である。第 1 の実施の形態では、例えば、第 1 のキーワード「日本」と第 2 のキーワード「面積」が入力されると、第 1 のキーワード「日本」に基づいて、第 1 のキーワードを、例えば「日本」、「アメリカ」、「ドイツ」という 3 つの第 3 のキーワードに増加させる。また、第 2 のキーワード「面積」に基づいて、第 2 のキーワードを、例えば「面積」、「人口」、「緯度」という 3 つの第 4 のキーワードに増加させる。そして、第 3 のキーワードと第 4 のキーワードとの組み合わせにより構成される、例えば「日本の面積は?」、「アメリカの人口は?」、「ドイツの緯度は?」・・・といった各質問に対する解答の候補を機械学習の手法を用いて求めて、解答として出力する。

【0073】

質問応答装置 1 は、入力されたキーワードを増加し、増加したキーワードにより構成される質問に対する解答を出力する装置である。質問応答装置 1 は、キーワード入力部 11

10

20

30

40

50

、キーワード増加部 1 2、質問作成部 1 3、解答候補抽出部 1 4、解答表出力部 1 5、キーワード抽出用データベース (DB) 1 6 を備える。また、図中、1 7 は後述する解答候補抽出部 1 4 による機械学習の結果 (学習結果) が蓄積されている学習データベース (DB) である。

【0074】

キーワード入力部 1 1 には、キーワードが入力される。例えば、第 1 のキーワード「日本」と第 2 のキーワード「面積」が入力される。キーワード増加部 1 2 は、後述するキーワード抽出技術を用いて、入力された各キーワードと同じ分野のキーワードをキーワード抽出用 DB 1 6 から抽出する。キーワードの抽出の結果、キーワードの総数が増加する。例えば、キーワード増加部 1 2 は、入力された第 1 のキーワードに基づいて、第 1 のキーワードの数より多い第 3 のキーワードを出力する。また、例えば、キーワード増加部 1 2 は、入力された第 2 のキーワードに基づいて、第 2 のキーワードの数より多い第 4 のキーワードを出力する。

10

【0075】

質問作成部 1 3 は、増加したキーワードである第 3 のキーワードと第 4 のキーワードとによって構成される質問を複数作成する。例えば、増加後の第 3 のキーワードの一つが「X」、第 4 のキーワードの一つが「Y」とすると、所属の格助詞「の」を用いて、「XのYは？」という質問を作成する。

【0076】

解答候補抽出部 1 4 は、後述する機械学習の手法によって、上記質問作成部 1 3 によって作成された質問に対する解答の候補である解答候補を抽出する。解答表出力部 1 5 は、抽出された各解答候補が質問と対応付けられた表を解答表として出力する。例えば、図 2 に示すような解答表を出力する。

20

【0077】

図 2 に示す解答表の例では、例えば、「日本の面積は？」という質問に対応する解答として、解答表のデータ項目「日本」に対応する行とデータ項目「面積」と対応する列とが交差する枠目に「A 1」(km²) が格納され、「アメリカの人口は？」という質問に対応する解答として、解答表のデータ項目「アメリカ」に対応する行とデータ項目「人口」と対応する列とが交差する枠目に「B 2」(万人) が格納される。

【0078】

本発明の実施の形態においては、抽出された解答候補を所定の単位 (例えば km²) に換算した表現を解答表に格納してもよく、また、抽出された解答候補についての単位のまま解答表に格納してもよい。

30

【0079】

もちろん、本発明において出力される解答表は、図 2 に示すものに限られるものではなく、例えば、「日本の面積は? A 1 (km²)」、「アメリカの人口は? B 2 (万人)」といった、各解答候補が矢印によって質問と対応付けられたデータが、解答表の各行のデータとして格納される形式の解答表を出力する構成を採ることもできる。

【0080】

キーワード抽出用 DB 1 6 は、一定量の文書データを格納したデータベースである。キーワード抽出用 DB 1 6 は、例えば、新聞、雑誌、Web データ (ネットワーク上のデータ) 等から抽出したデータ (一定量の文書データ) を格納している。学習 DB 1 7 には、後述する学習結果が蓄積されている。例えば、『質問「日本の首都は?」で答え「東京』という問題から抽出される素性の集合のときに、どのような解答 (「正解」または「不正解」) になりやすいかが、学習結果として蓄積されている。

40

【0081】

キーワード増加部 1 2 は、パターン抽出部 1 2 1 とキーワード抽出部 1 2 2 とを備える。パターン抽出部 1 2 1 は、キーワード入力部 1 1 に入力されたキーワードをキーワード抽出用 DB 1 6 で全文検索し、複数の入力キーワードの周辺に出現したパターンを抽出する。キーワード抽出部 1 2 2 は、パターン抽出部 1 2 1 で抽出したパターンをキーワード

50

抽出用DB16で全文検索し、該パターンによって抽出される表現をキーワードとして出力する。

【0082】

本発明の実施の形態においては、図1に示す構成から質問作成部13を省略し、解答候補抽出部14が、機械学習の手法を用いて、キーワード増加部12によって出力された第3のキーワードと第4のキーワードとによって構成される質問に対する解答候補を抽出し、出力する構成を採ってもよい。すなわち、解答候補抽出部14は、予め用意された問題と、その問題に対する解答の組の多数のセットを用いて、どういう問題のときにどういう解答になるかを学習し、その学習結果に基づいて、キーワード増加部12によって出力された第3のキーワードと第4のキーワードによって構成される質問に対する解答候補を抽出する構成を採ってもよい。

10

【0083】

以下に、キーワード増加部12によるキーワード抽出処理を説明する。パターン抽出部121は、入力された少数のキーワードをキーワード抽出用DB16で全文検索し、該少数のキーワードの周辺に出現したパターン c_i を抽出する。キーワード抽出部122は、抽出したパターン c_i をキーワード抽出用DB16で全文検索し、パターン c_i によって抽出される表現 exp を抽出すると同時に、抽出した表現 exp をScore(スコア; 評価値)の値の大きい順にソートしてキーワードとして出力する。

【0084】

本発明の実施の形態においては、キーワード抽出部122は、抽出した表現 exp について、Scoreの値が大きいものから順に所定の個数取り出してキーワードとして出力する構成を採ってもよい。また、キーワード抽出部122は、抽出した表現 exp について、Scoreの値が所定の閾値以上のものをキーワードとして出力する構成を採ってもよい。

20

【0085】

(パターンの例の説明)

以下に、パターン抽出部121が抽出するパターンについて、該パターンが国名Aである場合を例にとって説明する。

【0086】

・入力キーワード:

日本
中国
朝鮮
タイ
韓国

30

・抽出パターンの例(1): (両端とも利用、スピードは遅いが性能は良い)

日、A軍
人のA人女性
日本はAと
[A通信・
省。駐A大使な

40

・抽出パターンの例(2): (片方のみ利用、片方は平仮名文字、スピードは早い)

[..A国]。

【0087】

語。A
[..A国]側
[..A国]伝来
A語入力

ただし、[..A..]は、それ自体が国名Aにマッチすることを意味する。例えば[A国]だとそのマッチした用語の最後が国であることを意味する。

50

【0088】

(キーワード抽出の具体的な説明)

入力する少数のキーワードとして、例えば、評価データの代表形で毎日新聞での頻度の多い方から有名そうな用語を五つ選択するものとする。また、例えば、CD毎日新聞(コンパクトディスクに記録された毎日新聞)1991-2000年度版をキーワード抽出用DB16とする。抽出の手順例は以下のとおりである。

【0089】

(1) 少数の複数のキーワードをキーワード抽出用DB16で全文検索し、複数のキーワードの周辺に出現したパターンを c_i として抽出する(キーワードの周辺に出現するパターンがそのキーワードだけ(一個)の場合は抽出しない)。(周辺に出現するパターンの定義は適宜行なう)。周辺に出現するパターンとして例えば、キーワードの前後(左右)3文字列を用いる場合は、前後それぞれ文字が1個、2個、3個の場合があるので、1個のキーワードで9通りのパターンができることになる。また、キーワード(自分自身)を含めたパターンとすることもできる。

10

【0090】

(2) 次に抽出したパターン c_i をキーワード抽出用DB16で全文検索し、パターン c_i によって抽出される表現 exp を抽出する。

【0091】

(3) 抽出した表現 exp をScoreの値の大きい順にソートして、キーワードとして出力する。

20

【0092】

Scoreとして、以下のものがある。

【0093】

・手法1(決定リスト法)

手法1は、抽出した表現 exp のScoreとして、パターン c_i の中で p_i が最も大きかったパターンの p_i を使用する手法である。ここで、 p_i はパターン c_i で抽出される表現 exp での入力キーワードの割合(確からしさ、すなわち確信度となる)である。

【0094】

例えば、パターン c_1 についてキーワード抽出用DB16で全文検索した結果、 exp_1 、 exp_2 、 exp_3 、 exp_4 、 exp_5 までの5個の exp が抽出され、この5個の exp のうち、 $exp_1 \sim exp_3$ までの3個が入力キーワードであった場合、 p_1 は $3/5$ である。

30

【0095】

【数5】

$$\text{Score} = \max_i p_i \quad \text{式(7)}$$

【0096】

・手法2(ベイズ法)

手法2は、抽出した表現 exp のScoreとして、全てのパターン c_i の p_i を掛け合わせたものを使用する。

40

【0097】

【数6】

$$\text{Score} = \prod_i p_i \quad \text{式(8)}$$

【0098】

なお、実際には $p_i = 0$ の可能性が大きいため、本発明の実施の形態では、上記式(8)に代えて、以下の式(9)

50

$$\left(\frac{1 - p_i}{p_i + 1} \right) \quad \text{式(9)}$$

を利用する構成をとることもできる。ここで、 p_i は微小値の定数であり、例えば、0.0001を用いる。

【0099】

例えば、Scoreを計算しているexpが、パターン c_i についての検索処理によって取得できなかった場合は、 $p_i = 0$ として、上記の式(9)を用いて計算する。

【0100】

・手法3(類似度に基づく方法)

手法3は、抽出した表現expのScoreとして、抽出されたパターンの個数(総数)を用いる。つまり、多くのパターンで抽出されたものほどScoreを大きくする。

10

【0101】

【数7】

$$\text{Score} = \sum_i 1 \quad \text{式(10)}$$

【0102】

・手法4(下記研究(1)参照)

手法4は、抽出した表現expのScoreとして、 p_i の重みを加えた抽出されたパターンの個数を用いるものである。

20

【0103】

【数8】

$$\text{Score} = \sum_i (1 + 0.01p_i \log(f_i)) \quad \text{式(11)}$$

【0104】

ただし、 f_i はパターン c_i が出現した入力キーワードの個数である。

【0105】

研究(1): Ellen Riloff and Rosie Jones "Learning dictionaries for information extraction by multi-level bootstrapping" Proceedings of AAAI-99, (1999)。

30

【0106】

・手法5(下記文献(5)参照)

手法5は、抽出した表現expのScoreとして、少なくとも一つは確からしくなる値を用いるものである。

【0107】

【数9】

$$\text{Score} = 1 - \prod_i (1 - p_i) \quad \text{式(12)}$$

40

【0108】

上記式(12)は、確からしくない $(1 - p_i)$ を掛け合わせることで一つも確からしくないことになり、そして、これを1から引くと少なくとも一つは確からしくなる。

【0109】

文献(5): 村田真樹, 井佐原均 "同義テキストの照合に基づくパラフレーズに関する知識の自動獲得" 情報処理学会自然言語処理研究会 2001-NL-142, (2001)。

【0110】

上記手法1、2、4、5では、Scoreが同じときは、手法3のScoreでソート

50

し、手法3では手法5のScoreでソートする。

【0111】

図3は、パターンとしてキーワードの左と先頭のいずれかを含む1～3文字と右側のその組み合わせを用いて行ったキーワードの抽出結果に対して、予め用意した所定の種類の正解データを使って、適合率・再現率を求めた結果の一例を示す図である。ここで、正解データとしては、例えば、図4に示すようなデータ例を用意する(図4は、国名データの例を示しており、国名を国ごとに行に分けて格納し、行頭を代表形としてそれ以外は代表形の異表記として同じ行に格納している)。図4に示すデータ形式と同様のデータ形式を持つ正解データを、例えば、国名データの他に、衛星、祝日、太陽系惑星、世界遺産等に関するデータのように、多種類用意する。

10

【0112】

図3において、APは、情報検索(下記文献(6)参照)で用いるaverage precisionの平均であり、正解記事を上位から取った時に求めた適合率の平均である。本願の内容の場合は、正解キーワード分を上位から取った時に求めた適合率の平均(ただし、入力キーワードは正解キーワードから除く)である。

【0113】

文献(6): 村田真樹, 馬青, 内元清貴, 小作浩美, 内山将夫, 井佐原均 "位置情報と分野情報を用いた情報検索" 言語処理学会誌, Vol.7, No.2, (2000)。

【0114】

RPは、r-precisionの平均であり、正解記事数分だけを検索した時に正解の記事が含まれている割合である。本願の内容の場合は、正解キーワード分だけを抽出した時に正解キーワードが含まれている割合である。なお、適合率は正解率と同じであり、正解キーワードが含まれる割合のことである。TPは、上位5個での精度の平均である。

20

【0115】

(制約に基づく抽出方法の説明)

(a) 字種とKRを利用する方法

図3に示す例で、抽出方法には、さらに字種とKRを利用する方法を用いた。ここで、字種とは、漢字、カタカナ、ひらがな、記号、数字などであり、例えば英語だと、アルファベット、数字、記号、単語の先頭が大文字かどうかなどである。

【0116】

字種を利用する方法では、入力した少数(この例では5個)のキーワードになかった字種を含む表現を抽出しない方法である。例えば、入力した5個のキーワードにひらがなが無かった場合は、ひらがなを含む表現を抽出しないようにするものである。

30

【0117】

KRを利用する方法では、 p_i を $p_i * f_i / n_i$ に置き換えた方法である。この方法の利点は、 p_i が同じでも f_i / n_i の値により確信度を変えることができるものである。ただし、 n_i は入力キーワードの個数で、手法3のときはKRの場合は1を f_i に置き換えた。なお、評価では抽出した結果でキーワードの異表記は除いた。また、字種による方法以外にも次のような方法もある。

【0118】

(b) 品詞に基づく方法

品詞に基づく方法では、例えば、入力表現に名詞しかない場合は出力時に名詞以外の表現を省く、また、入力表現に形容詞しかない場合は出力時に形容詞以外の表現を省くというものである。さらに、表現が複数の単語で構成されている場合は、末尾の単語(形態素)の品詞の情報を使うようにすることができる。

40

【0119】

(例による説明1)

入力キーワードとして次のものであった場合、

「楽しい」「哀しい」「嬉しい」「とても嬉しい」「とても哀しい」

抽出物として次のものが得られる場合、

50

「とても」「新しい」「美しい」「とても美しい」「とても難しい」

上記抽出物の表現中の末尾の単語の品詞を推定し、上記入力キーワードでは、末尾の単語の品詞は「形容詞」しかないので、抽出物の中で、末尾の単語の品詞が「形容詞」でない、副詞（「とても」）を除いて出力するようにする。

【0120】

（例による説明2）

入力キーワードとして次のものであった場合、

「楽しい」「歓喜」「悲痛」「悲しい」

上記入力キーワードでは、「形容詞」と「名詞」のように複数種類があった場合は、それらの品詞は出力し、それらの品詞以外の表現は出力しないようにする。

10

【0121】

なお、前述のような末尾の単語（形態素）の品詞の推定等の品詞情報を得るためには、次のような形態素解析システム（形態素解析手段）が必要になる。

【0122】

・形態素解析システムの説明

日本語を単語に分割するために、キーワード抽出部122で形態素解析システムを利用することが必要になる。ここではChaSenについて説明する（奈良先端大で開発されている形態素解析システム茶筌。<http://chasen.aist-nara.ac.jp/index.html.jp> で公開されている）。

【0123】

これは、日本語文を分割し、さらに、各単語の品詞も推定してくれる。例えば、「学校へ行く」を入力すると以下の結果を得ることができる。

20

【0124】

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

E O S

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

【0125】

（c）共通部分文字列に基づく方法

例えば、入力表現がすべて同じ「しい」という共通末尾表現を持っている場合、出力時に「しい」を持たない表現を省くものである。なお、これは末尾だけでなく、先頭の文字列でも同様にできる。

30

【0126】

（例による説明）

入力キーワードとして次のものであった場合、

「悲しい」「楽しい」「嬉しい」

抽出されるものが次の場合、

「歓喜」「悲痛」「美しい」「新しい」

上記入力キーワードの共通部分文字列が「しい」なので、「しい」を持たない「歓喜」と「悲痛」を削除して出力するものである。

40

【0127】

（d）ユーザによる制約の指定

上記では、入力表現から自動で制約を得る方法を説明したが、この制約はユーザにさせることもできる。例えば、ユーザが「漢字のみ」というオプションを選択すると出力では漢字以外の字種を用いた表現を出力しないことができる。また、ユーザが末尾は「しい」というオプションを選択すると出力では「しい」を末尾を持たない表現を出力しないようにすることができる。さらに、ユーザが品詞は名詞というオプションを選択すると出力では名詞以外の表現を出力しないようにする。

50

【 0 1 2 8 】

次に、質問作成部 1 3 が作成した質問、または、質問応答装置 1 が質問作成部 1 3 を備えない構成を採るときはキーワード増加部 1 2 によって出力された出力キーワードによって構成される質問（「XのYは？」）に対する解答候補を抽出する処理について説明する。解答候補抽出部 1 4 は、機械学習の手法を用いて解答候補を抽出する。

（機械学習の手法）

機械学習の手法は、問題 - 解答の組のセットを多く用意し、それで学習を行ない、どういう問題のときにどういう解答になるかを学習し、その学習結果を利用して、新しい問題のときも解答を推測できるようにする方法である（例えば、下記の文献（7）参照）。

【 0 1 2 9 】

文献（7）：村田真樹，機械学習に基づく言語処理，龍谷大学理工学部．招待講演．2004．<http://www2.nict.go.jp/jt/a132/members/murata/ps/rk1-siryoku.pdf>

どういう問題のときに、という、問題の状況を機械に伝える際に、素性（解析に用いる情報で問題を構成する各要素）というものがことになる。問題を素性によって表現するのである。例えば、日本語文末表現の時制の推定の問題において、問題：「彼が話す。」 - - 解答「現在」が与えられた場合に、素性の一例は、「彼が話す。」「が話す。」「話す。」「す」「。」となる。

【 0 1 3 0 】

すなわち、機械学習の手法は、素性の集合 - 解答の組のセットを多く用意し、それで学習を行ない、どういう素性の集合のときにどういう解答になるかを学習し、その学習結果を利用して、新しい問題のときもその問題から素性の集合を取り出し、その素性の場合の解答を推測する方法である。

【 0 1 3 1 】

まず、機械学習の手法一般についての説明をする。機械学習の手法としては、一般に、k近傍法、シンプルベイズ法、決定リスト法、最大エントロピー法、サポートベクトルマシン法などの手法を用いる。

【 0 1 3 2 】

k近傍法は、最も類似する一つの事例のかわりに、最も類似するk個の事例を用いて、このk個の事例での多数決によって分類先（解）を求める手法である。kは、あらかじめ定める整数の数字であって、一般的に、1から9の間の奇数を用いる。

【 0 1 3 3 】

シンプルベイズ法は、ベイズの定理にもとづいて各分類になる確率を推定し、その確率値が最も大きい分類を求める分類先とする方法である。

【 0 1 3 4 】

シンプルベイズ法において、文脈bで分類aを出力する確率は、以下の式（13）で与えられる。

【 0 1 3 5 】

【数 1 0】

$$p(a|b) = \frac{p(a)}{p(b)} p(b|a) \quad \text{式(13)}$$

$$\cong \frac{\tilde{p}(a)}{p(b)} \prod_i \tilde{p}(f_i|a) \quad \text{式(14)}$$

【 0 1 3 6 】

ただし、ここで文脈bは、あらかじめ設定しておいた素性 f_j （ $F, 1 \leq j \leq k$ ）の

10

20

30

40

50

集合である。 $p(b)$ は、文脈 b の出現確率である。ここで、分類 a に非依存であって定数のために計算しない。 $P(a)$ (ここで P は p の上部にチルダ) と $P(f_i | a)$ は、それぞれ教師データから推定された確率であって、分類 a の出現確率、分類 a のときに素性 f_i を持つ確率を意味する。 $P(f_i | a)$ として最尤推定を行って求めた値を用いると、しばしば値がゼロとなり、式(14)の値がゼロで分類先を決定することが困難な場合が生じる。そのため、スムージングを行う。ここでは、以下の式(15)を用いてスムージングを行ったものを用いる。

【0137】

【数11】

$$p(f_i | a) = \frac{\text{freq}(f_i, a) + 0.01 * \text{freq}(a)}{\text{freq}(a) + 0.01 * \text{freq}(a)} \quad \text{式(15)}$$

10

【0138】

ただし、 $\text{freq}(f_i, a)$ は、素性 f_i を持ちかつ分類が a である事例の個数、 $\text{freq}(a)$ は、分類が a である事例の個数を意味する。

【0139】

決定リスト法は、素性と分類先の組とを規則とし、それらをあらかじめ定めた優先順序でリストに蓄えおき、検出する対象となる入力を与えられたときに、リストで優先順位の高いところから入力のデータと規則の素性とを比較し、素性が一致した規則の分類先をその入力の分類先とする方法である。

20

【0140】

決定リスト方法では、あらかじめ設定しておいた素性 f_j ($F, 1 \leq j \leq k$) のうち、いずれか一つの素性のみを文脈として各分類の確率値を求める。ある文脈 b で分類 a を出力する確率は以下の式によって与えられる。

【0141】

$$p(a | b) = p(a | f_{\max}) \quad \text{式(16)}$$

ただし、 f_{\max} は以下の式によって与えられる。

【0142】

【数12】

$$f_{\max} = \arg \max_{f_j \in F} \max_{a_i \in A} \tilde{p}(a_i | f_j) \quad \text{式(17)}$$

30

【0143】

また、 $P(a_i | f_j)$ (ここで P は p の上部にチルダ) は、素性 f_j を文脈に持つ場合の分類 a_i の出現の割合である。

【0144】

最大エントロピー法は、あらかじめ設定しておいた素性 f_j ($1 \leq j \leq k$) の集合を F とするとき、以下所定の条件式(式(18))を満足しながらエントロピーを意味する式(19)を最大にするときの確率分布 $p(a, b)$ を求め、その確率分布にしたがって求める各分類の確率のうち、最も大きい確率値を持つ分類を求める分類先とする方法である。

40

【0145】

【数 1 3】

$$\sum_{a \in A, b \in B} p(a, b) g_j(a, b) = \sum_{a \in A, b \in B} \tilde{p}(a, b) g_j(a, b) \quad \text{式(18)}$$

for $\forall f_j (1 \leq j \leq k)$

10

$$H(p) = - \sum_{a \in A, b \in B} p(a, b) \log(p(a, b)) \quad \text{式(19)}$$

【0 1 4 6】

ただし、A、Bは分類と文脈の集合を意味し、 $g_j(a, b)$ は文脈bに素性 f_j があつて、なおかつ分類がaの場合1となり、それ以外で0となる関数を意味する。また、 $P(a_i | f_j)$ (ここでPはpの上部にチルダ)は、既知データでの(a, b)の出現の割合を意味する。

20

【0 1 4 7】

式(18)は、確率pと出力と素性の組の出現を意味する関数gをかけることで出力と素性の組の頻度の期待値を求めることになっており、右辺の既知データにおける期待値と、左辺の求める確率分布に基づいて計算される期待値が等しいことを制約として、エントロピー最大化(確率分布の平滑化)を行なつて、出力と文脈の確率分布を求めるものとなっている。最大エントロピー法の詳細については、以下の文献(8)および文献(9)に記載されている。

【0 1 4 8】

文献(8) : Eric Sven Ristad, Maximum Entropy Modeling for Natural Language, (ACL/EACL Tutorial Program, Madrid, 1997)

30

文献(9) : Eric Sven Ristad, Maximum Entropy Modeling Toolkit, Release 1.6beta, (<http://www.mnemonic.com/software/memt>, 1998)

サポートベクトルマシン法は、空間を超平面で分割することにより、二つの分類からなるデータを分類する手法である。

【0 1 4 9】

図27にサポートベクトルマシン法のマージン最大化の概念を示す。図27において、白丸は正例、黒丸は負例を意味し、実線は空間を分割する超平面を意味し、破線はマージン領域の境界を表す面を意味する。図27(A)は、正例と負例の間隔が狭い場合(スモールマージン)の概念図、図27(B)は、正例と負例の間隔が広い場合(ラージマージン)の概念図である。

40

【0 1 5 0】

このとき、二つの分類が正例と負例からなるものとする、学習データにおける正例と負例の間隔(マージン)が大きいものほどオープンデータで誤った分類をする可能性が低いと考えられ、図27(B)に示すように、このマージンを最大にする超平面を求めそれを用いて分類を行なう。

【0 1 5 1】

基本的には上記のとおりであるが、通常、学習データにおいてマージンの内部領域に少数の事例が含まれてもよいとする手法の拡張や、超平面の線形の部分を非線形にする拡張(カーネル関数の導入)がなされたものが用いられる。

50

【 0 1 5 2 】

この拡張された方法は、以下の識別関数を用いて分類することと等価であり、その識別関数の出力値が正か負かによって二つの分類を判別することができる。

【 0 1 5 3 】

【数 1 4】

$$f(x) = \operatorname{sgn}\left(\sum_{i=1}^l \alpha_i y_i K(x_i, x) + b\right) \quad \text{式(20)}$$

10

$$b = -\frac{\max_{i, y_i=-1} b_i + \min_{i, y_i=1} b_i}{2}$$

$$b_i = \sum_{j=1}^l \alpha_j y_j K(x_j, x_i)$$

【 0 1 5 4 】

ただし、 x は識別したい事例の文脈（素性の集合）を、 x_i と y_j ($i = 1, \dots, l$, $y_j \in \{1, -1\}$) は学習データの文脈と分類先を意味し、関数 sgn は、

$$\operatorname{sgn}(x) = \begin{cases} 1 & (x \geq 0) \\ -1 & (\text{otherwise}) \end{cases}$$

であり、また、各 α_i は式(22)と式(23)の制約のもと式(21)を最大にする場合のものである。

【 0 1 5 5 】

【数 1 5】

$$L(\alpha) = \sum_{i=1}^l \alpha_i - \frac{1}{2} \sum_{i,j=1}^l \alpha_i \alpha_j y_i y_j K(x_i, x_j) \quad \text{式(21)}$$

30

$$0 \leq \alpha_i \leq C \quad (i = 1, \dots, l) \quad \text{式(22)}$$

$$\sum_{j=1}^l \alpha_j y_j = 0 \quad \text{式(23)}$$

【 0 1 5 6 】

また、関数 K はカーネル関数と呼ばれ、様々なものが用いられるが、本形態では以下の多項式のものを用いる。

【 0 1 5 7 】

$$K(x, y) = (x \cdot y + 1)^d \quad \text{式(24)}$$

C 、 d は実験的に設定される定数である。後述する具体例では C はすべての処理を通して 1 に固定した。また、 d は、1 と 2 の二種類を試している。ここで、 $\alpha_i > 0$ となる x_i は、サポートベクトルと呼ばれ、通常、式(20)の和をとっている部分は、この事例のみを用いて計算される。つまり、実際の解析には学習データのうちサポートベクトルと呼ばれる事例のみしか用いられない。

【 0 1 5 8 】

40

50

なお、拡張されたサポートベクトルマシン法の詳細については、以下の文献(10)および文献(11)に記載されている。

【0159】

文献(10) : Nello Cristianini and John Shawe-Taylor, An Introduction to Support Vector Machines and other kernel-based learning methods, (Cambridge University Press, 2000)

文献(11) : Taku Kudoh, Tinsvm: Support Vector machines, (<http://cl.aist-nara.ac.jp/taku-ku//software/Tiny SVM/index.html>, 2000)

サポートベクトルマシン法は、分類の数が2個のデータを扱うものである。したがって、分類の数が3個以上の事例を扱う場合には、通常、これにペアワイズ法またはワンVS

10

レスト法などの手法を組み合わせて用いることになる。

【0160】

ペアワイズ法は、 n 個の分類を持つデータの場合に、異なる二つの分類先のあらゆるペア($n(n-1)/2$ 個)を生成し、各ペアごとにどちらがよいかを二値分類器、すなわちサポートベクトルマシン法処理モジュールで求めて、最終的に、 $n(n-1)/2$ 個の二値分類による分類先の多数決によって、分類先を求める方法である。

【0161】

ワンVSレスト法は、例えば、 a 、 b 、 c という三つの分類先があるときは、分類先 a とその他、分類先 b とその他、分類先 c とその他、という三つの組を生成し、それぞれの組についてサポートベクトルマシン法で学習処理する。そして、学習結果による推定処理において、その三つの組のサポートベクトルマシンの学習結果を利用する。推定すべき二項関係の候補が、その三つのサポートベクトルマシンではどのように推定されるかを見て、その三つのサポートベクトルマシンのうち、その他でないほうの分類先であって、かつサポートベクトルマシンの分離平面から最も離れた場合のものの分類先を求める解とする方法である。例えば、ある候補が、「分類先 a とその他」の組の学習処理で作成したサポートベクトルマシンにおいて分離平面から最も離れた場合には、その候補の分類先は a と推定する。

20

【0162】

次に、本発明の実施の形態における具体的な解答候補の抽出手法を説明する。

(解答候補の抽出手法1)

30

<問題の構成>

予め、

問題『質問「 X_1 の Y_1 は？」で答え「 Z_1 」』 - - - 解答「正解」

問題『質問「 X_2 の Y_2 は？」で答え「 Z_2 」』 - - - 解答「正解」

問題『質問「 X_3 の Y_3 は？」で答え「 Z_3 」』 - - - 解答「不正解」

という、問題と解答の対を多数作成する。

【0163】

また、例えば、上記の問題を表現する、以下のような素性を用意する。

- ・ X_i , Y_i , Z_i の単語自体
- ・ X_i , Y_i , Z_i の単語の意味クラス
- ・ X_i , Y_i で検索した記事数
- ・ X_i , Y_i で検索した記事に Z_i が存在する記事数
- ・ X_i , Y_i が近接して(ある単語数の範囲内に)出現した記事数
- ・ X_i , Y_i , Z_i が近接して(ある単語数の範囲内に)出現した記事数
- ・ X_i , Y_i で検索した記事に最も多く出現した単語と Z_i が一致するかどうか
- ・ X_i , Y_i で検索した記事に j 番目に多く出現した単語と Z_i が一致するかどうか
- ・ X_i , Y_i をキーワードとして、例えば、解答候補抽出部14が、新聞記事データ・百科事典データなどの文書データ群(図示を省略)から解答の書いてありそうな記事群を取り出し、その取り出した記事群の言語表現を解答の候補として取り出し、取り出された解答の候補を、優先順序(例えば、 $Score_{near}(c)$)で並び替えた場合に、その順序

40

50

の最も高い候補と Z_i が一致したかどうか、また、その順序の j 番目の候補と Z_i が一致したかどうか

上記の処理によって、素性の集合と解答の組の多数のセットが用意される。

【0164】

ここで、優先順序として用いる $Score_{near}(c)$ については、前述の非特許文献1に記述されており、解答の候補とキーワードの近さに基づく得点を示している。

【0165】

次に、意味クラスを説明する。一般に、各単語がどういう意味クラスを持つかを記述した表があり、その表を使えば、単語の意味クラスを求めることができる。例えば分類語彙表がある。分類語彙表では単語は分類番号と呼ばれる10桁の数字で表現され、この数字の良く似ている単語ほど良く似た単語となる。例えば、この数字の最初の3桁や5桁を単語の意味クラスとして利用する。例えば、「村人」の分類番号は1230102050であり、これは123(人種、国民、社会階層などの意味クラス)、12301(国民、住民などの意味クラス)に属する単語であることが示される。

【0166】

問題構成と素性の定義をすれば、あとは機械学習の手法で扱える。すなわち、解答候補抽出部14は、用意された素性と解答の組の多数のセットを用いて、どういう素性の集合のときにどういう解答になるかを学習し、その学習結果を利用して、新たな問題についての素性の集合の場合に推測される解答を、解答候補として抽出する。

【0167】

<問題や素性の具体例>

問題の具体例：

問題『質問「日本の首都は？」で答え「東京」』 - - - 解答「正解」
 問題『質問「日本の首都は？」で答え「大阪」』 - - - 解答「不正解」
 問題『質問「日本の首都は？」で答え「パン」』 - - - 解答「不正解」

素性の具体例：

問題『質問「日本の首都は？」で答え「東京」』 - - - 解答「正解」の場合

- ・ X_i の単語自体：日本
- ・ Y_i の単語自体：首都
- ・ Z_i の単語自体：東京

- ・ X_i の意味クラス：12590(地名のクラス)
- ・ Y_i の意味クラス：12540(都市集落のクラス)
- ・ Z_i の意味クラス：12590(地名のクラス)

(意味クラスとして分類語彙表の最初の5桁を利用する。)

- ・ X_i, Y_i で検索した記事数：日本と首都を含む記事数。例えば1000
- ・ X_i, Y_i で検索した記事に Z_i が存在する記事数：日本と首都と東京を含む記事数。

例えば100

・ X_i, Y_i が近接して(ある単語数の範囲内に)出現した記事数：例えば、日本と首都が10単語以内にある記事数。例えば500

・ X_i, Y_i, Z_i が近接して(ある単語数の範囲内に)出現した記事数：例えば、日本と首都と東京が10単語以内にある記事数。例えば50

・ X_i, Y_i で検索した記事に最も多く出現した単語と Z_i が一致するかどうか：例えば、ここでは特に単語は名詞にしばり、名詞としては、「こと」が最も頻度が多かったとすると、「こと」と「東京」が一致しないのでこの素性は「いいえ」となる。

・ X_i, Y_i で検索した記事に j 番目に多く出現した単語と Z_i が一致するかどうか：例えば、ここでは特に単語は名詞にしばり、名詞としては、「東京」が二番目に頻度が多かったとすると、 $j=2$ の場合の素性は「はい」となる。

・ X_i, Y_i をキーワードとして、例えば、解答候補抽出部14が、新聞記事データ・百科事典データなどの文書データ群から解答の書いてありそうな記事群を取り出し、その取り出した記事群の言語表現を解答の候補を取り出し、取り出された解答の候補を、優先順

10

20

30

40

50

序（例えば、 $Score_{near}(c)$ ）で並び替えた場合に、その順序の最も高い候補と Z_i が一致したかどうか、また、その順序の j 番目の候補と Z_i が一致したかどうか

例えば、1 番目の候補が「こと」、2 番目の候補が「東京」の場合は、「その順序の最も高い候補と Z_i が一致したかどうか」は「いいえ」になり、「その順序の j 番目の候補と Z_i が一致したかどうか」は、 $j = 2$ のとき「はい」になる。

【0168】

より多くの事例で学習すると、例えば、解答候補抽出部 14 は、

Y_i の単語自体：首都

Z_i の意味クラス：12590（地名のクラス）

で、「その順序の最も高い候補と Z_i が一致したかどうか」は「はい」

10

または「その順序の j 番目の候補と Z_i が一致したかどうか」の $j = 2$ のときが「はい」になれば、

解答「正解」

となるように学習し、

Y_i の単語自体：首都

Z_i の意味クラス：12590（地名のクラス）以外なら、

解答「不正解」

Y_i の単語自体：首都

「その順序の最も高い候補と Z_i が一致したかどうか」は「いいえ」かつ、

「その順序の j 番目の候補と Z_i が一致したかどうか」の $j = 2$ から 10 全てで「いいえ」ならば、

20

解答「不正解」

といったことを学習する。

【0169】

学習結果は、解答候補抽出部 14 によって、学習 DB 17 中に蓄積される。そして、解答候補抽出部 14 は、学習 DB 17 中の学習結果情報を使い、例えば、新しい問題（すなわち、キーワード増加部 12 によって出力された出力キーワードによって構成される質問）：

質問「フランスの首都は？」 - 答え「パリ」については、

Y_i の単語自体：首都

30

Z_i の意味クラス：12590（地名のクラス）で、

「その順序の最も高い候補と Z_i が一致したかどうか」は「はい」または、

「その順序の j 番目の候補と Z_i が一致したかどうか」の $j = 2$ のときが「はい」

なので、「正解」と判断する。

【0170】

また、新しい問題：

質問「フランスの首都は？」 - 答え「信号」については、上記学習結果を用いて、

Y_i の単語自体：首都

Z_i の意味クラス：12590（地名のクラス）以外

なので、「不正解」と判断する。

40

【0171】

ここで、機械学習の手法によって、解答を求めるだけでなく、その解答がどのくらい正解になりやすいかの度合い、どのくらい不正解になりやすいかの度合いも同時に求めることができる。

【0172】

すなわち、解答候補抽出部 14 は、用意した素性の集合と解答の組の多数のセットを用いて、まず、どのような素性の集合のときにどのような解答（正解または不正解）となるかということを経験学習し、どのような素性の集合のときにどのような解答となるかということを示す情報を、学習結果情報として学習 DB 17 に格納する。そして、解答候補抽出部 14 は、新たな問題（キーワード増加部 12 によって出力された出力キーワードによ

50

って構成される質問)を用いて、新たに解答候補抽出部14が作成した問題)から素性の集合を抽出し、抽出された素性の集合の場合にどのような解答になりやすいか、すなわち、「正解となりやすい」かの度合いを、学習DB17に格納された学習結果情報に基づいて求める。

【0173】

そして、解答候補抽出部14は、例えば、「正解となりやすい」かの度合いが最も大きいときの、問題(質問-答えの対)における、「答え」を、解答候補として解答表出力部15に対して出力する。解答表出力部15は、各解答候補が質問と対応付けられた表を解答表として出力する。本発明の実施の形態においては、例えば、「正解となりやすい」かの度合いの大きい順に所定の個数の問題を選択し、選択した問題における「答え」を解答候補としてもよい。また、例えば、「正解となりやすい」かの度合いが所定の閾値以上の問題を選択し、選択した問題における「答え」を解答候補としてもよい。また、例えば、「正解となりやすい」かの度合いが最も大きい問題についての当該度合いの所定の割合(例えば、90%等)を閾値とし、「正解となりやすい」かの度合いがこの閾値以上の問題を選択し、選択した問題における「答え」を解答候補としてもよい。

10

【0174】

上記の、「正解となりやすい」かの度合いの求め方は、解答候補抽出部14が機械学習の手法として用いる様々な方法によって異なる。

【0175】

例えば、本発明の実施の形態において、解答候補抽出部14が、機械学習の手法としてk近傍法を用いる場合、解答候補抽出部14は、上記用意した素性の集合と解答の組の多数のセットを用いて、素性の集合のうち重複する素性の割合(同じ素性をいくつ持っているかの割合)に基づく問題同士の類似度を定義する。そして、解答候補抽出部14は、定義した類似度と問題(と解答)とを、学習結果情報として、学習DB17に格納しておく。

20

【0176】

そして、解答候補抽出部14は、質問作成部13によって作成された質問、または、質問応答装置1が質問作成部13を備えない構成を採るときはキーワード増加部12によって出力された出力キーワードによって構成される質問を用いて新たに問題を作成すると、学習DB17内に格納された類似度と問題を参照して、新たに作成された問題との類似度が高い順にk個の問題(と解答)を学習DB17に格納された問題(と解答)から選択し、選択したk個の問題での多数決によって決まった分類先(正解または不正解)を、新たに作成された問題に対する解答とする。k近傍法を用いる場合、「正解となりやすい」かの度合いは、上記選択されたk個の問題での多数決の票数、すなわち、「正解」という分類が獲得した票数となる。

30

【0177】

解答候補抽出部14は、「正解となりやすい」かの度合いが最も大きいときの、問題(質問-答えの対)における、「答え」を、解答候補として解答表出力部15に対して出力する。本発明の実施の形態においては、例えば、「正解となりやすい」かの度合いの大きい順に所定の個数の問題を選択し、選択した問題における「答え」を解答候補としてもよい。また、例えば、「正解となりやすい」かの度合いが所定の閾値以上の問題を選択し、選択した問題における「答え」を解答候補としてもよい。また、例えば、「正解となりやすい」かの度合いが最も大きい問題についての当該度合いの所定の割合(例えば、90%等)を閾値とし、「正解となりやすい」かの度合いがこの閾値以上の問題を選択し、選択した問題における「答え」を解答候補としてもよい。

40

【0178】

また、例えば、本発明の実施の形態において、解答候補抽出部14が、機械学習の手法としてシンプルベイズ法を用いる場合、例えば、解答候補抽出部14は、上記用意した素性の集合と解答の組の多数のセットを学習結果情報として学習DB17に格納しておく。

【0179】

50

解答候補抽出部 14 は、質問作成部 13 によって作成された質問、または、質問応答装置 1 が質問作成部 13 を備えない構成を採るときはキーワード増加部 12 によって出力された出力キーワードによって構成される質問を用いて新たに問題を作成すると、新たに作成した問題から素性の集合を抽出する。そして、解答候補抽出部 14 は、学習 DB 17 内に格納された解答と素性の集合とのセットをもとに、ベイズの定理に基づいて、新たに作成された問題から抽出した素性の集合の場合の各分類になる確率を算出して、その確率の値が最も大きい分類を、その問題に対する解答とする。シンプルベイズ法を用いる場合、「正解となりやすい」かの度合いは、「正解」という分類になる確率となる。

【0180】

解答候補抽出部 14 は、「正解となりやすい」かの度合いが最も大きいときの、問題（質問 - 答えの対）における、「答え」を、解答候補として解答表出力部 15 に対して出力する。本発明の実施の形態においては、例えば、「正解となりやすい」かの度合いの大きい順に所定の個数の問題を選択し、選択した問題における「答え」を解答候補としてもよい。また、例えば、「正解となりやすい」かの度合いが所定の閾値以上の問題を選択し、選択した問題における「答え」を解答候補としてもよい。また、例えば、「正解となりやすい」かの度合いが最も大きい問題についての当該度合いの所定の割合（例えば、90%等）を閾値とし、「正解となりやすい」かの度合いがこの閾値以上の問題を選択し、選択した問題における「答え」を解答候補としてもよい。

10

【0181】

また、例えば、本発明の実施の形態において、解答候補抽出部 14 が、機械学習の手法として決定リスト法を用いる場合、例えば、解答候補抽出部 14 は、予め用意した問題についての素性と分類先との規則を所定の優先順序で並べたリストを学習 DB 17 内に格納する。解答候補抽出部 14 は、質問作成部 13 によって作成された質問、または、質問応答装置 1 が質問作成部 13 を備えない構成を採るときはキーワード増加部 12 によって出力された出力キーワードによって構成される質問を用いて新たに問題を作成すると、新たに作成した問題から素性の集合を抽出する。

20

【0182】

そして、解答候補抽出部 14 は、学習 DB 17 内に格納されたリストの優先順位の高い順に、上記新たに作成した問題から抽出された素性と規則の素性とを比較し、素性が一致した規則の分類先をその問題に対する解答とする。決定リスト法を用いる場合、「正解となりやすい」かの度合いは、所定の優先順位またはそれに相当する数値、尺度となる。

30

【0183】

解答候補抽出部 14 は、「正解となりやすい」かの度合いが最も大きいときの、問題（質問 - 答えの対）における、「答え」を、解答候補として解答表出力部 15 に対して出力する。本発明の実施の形態においては、例えば、「正解となりやすい」かの度合いの大きい順に所定の個数の問題を選択し、選択した問題における「答え」を解答候補としてもよい。また、例えば、「正解となりやすい」かの度合いが所定の閾値以上の問題を選択し、選択した問題における「答え」を解答候補としてもよい。また、例えば、「正解となりやすい」かの度合いが最も大きい問題についての当該度合いの所定の割合（例えば、90%等）を閾値とし、「正解となりやすい」かの度合いがこの閾値以上の問題を選択し、選択した問題における「答え」を解答候補としてもよい。

40

【0184】

また、例えば、本発明の実施の形態において、解答候補抽出部 14 が、機械学習の手法として最大エントロピー法を用いる場合、例えば、解答候補抽出部 14 は、予め用意した問題の解答となりうる分類を特定し、所定の条件式を満足しかつエントロピーを示す式を最大にするときの素性の集合と解答となりうる分類の二項からなる確率分布を求めて、学習 DB 17 内に格納する。

【0185】

そして、解答候補抽出部 14 は、質問作成部 13 によって作成された質問、または、質問応答装置 1 が質問作成部 13 を備えない構成を採るときはキーワード増加部 12 によ

50

て出力された出力キーワードによって構成される質問を用いて新たに問題を作成すると、学習DB17内に格納された確率分布を利用して、新たな問題の素性の集合についてその解答となりうる分類の確率を求めて、最も大きい確率値を持つ解答となりうる分類を特定し、その特定した分類をその問題に対する解答とする。

【0186】

すなわち、最大エントロピー法を用いる場合、「正解となりやすい」かの度合いは、「正解」という分類になる確率となる。

【0187】

解答候補抽出部14は、「正解となりやすい」かの度合いが最も大きいときの、問題（質問 - 答えの対）における、「答え」を、解答候補として解答表出力部15に対して出力する。本発明の実施の形態においては、例えば、「正解となりやすい」かの度合いの大きい順に所定の個数の問題を選択し、選択した問題における「答え」を解答候補としてもよい。また、例えば、「正解となりやすい」かの度合いが所定の閾値以上の問題を選択し、選択した問題における「答え」を解答候補としてもよい。また、例えば、「正解となりやすい」かの度合いが最も大きい問題についての当該度合いの所定の割合（例えば、90%等）を閾値とし、「正解となりやすい」かの度合いがこの閾値以上の問題を選択し、選択した問題における「答え」を解答候補としてもよい。

10

【0188】

また、例えば、本発明の実施の形態において、解答候補抽出部14が、機械学習の手法としてサポートベクトルマシン法を用いる場合、例えば、解答候補抽出部14は、予め用意した問題の解答となりうる分類を特定し、分類を正例と負例に分割して、カーネル関数を用いた所定の実行関数に従って問題の素性の集合を次元とする空間上で、その問題の正例と負例の間隔を最大にし、かつ正例と負例を超平面で分割する超平面を求めて学習DB17内に格納する。

20

【0189】

そして、解答候補抽出部14は、質問作成部13によって作成された質問、または、質問応答装置1が質問作成部13を備えない構成を採るときはキーワード増加部12によって出力された出力キーワードによって構成される質問を用いて新たに問題を作成すると、学習DB17内の超平面を利用して、新たな問題の素性の集合が超平面で分割された空間において正例側か負例側のどちらにあるかを特定し、その特定された結果に基づいて定まる分類を、その問題に対する解答とする。

30

【0190】

すなわち、サポートベクトルマシン法を用いる場合、「正解となりやすい」かの度合いは、分離平面からの正例の空間への距離の大きさとなる。より詳しくは、解答が正解である問題を正例、解答が不正解である問題を負例とする場合に、分離平面に対して正例側の空間に位置する問題が、解答が正解である問題と判断され、分離平面からの距離が大きい問題ほど「正解となりやすい」かの度合いが大きくなる。

【0191】

解答候補抽出部14は、「正解となりやすい」かの度合いが最も大きいときの、問題（質問 - 答えの対）における、「答え」を、解答候補として解答表出力部15に対して出力する。本発明の実施の形態においては、例えば、「正解となりやすい」かの度合いの大きい順に所定の個数の問題を選択し、選択した問題における「答え」を解答候補としてもよい。また、例えば、「正解となりやすい」かの度合いが所定の閾値以上の問題を選択し、選択した問題における「答え」を解答候補としてもよい。また、例えば、「正解となりやすい」かの度合いが最も大きい問題についての当該度合いの所定の割合（例えば、90%等）を閾値とし、「正解となりやすい」かの度合いがこの閾値以上の問題を選択し、選択した問題における「答え」を解答候補としてもよい。

40

【0192】

<具体例>

例えば、質問作成部13によって作成された質問、または、質問応答装置1が質問作成

50

部 1 3 を備えない構成を採るときはキーワード増加部 1 2 によって出力された出力キーワードによって構成される質問「フランスの首都は？」を解く場合について説明する。

【 0 1 9 3 】

まず、解答候補抽出部 1 4 が、出力キーワード「フランス」、「首都」を含む文書を、新聞記事データ・百科事典データなどの文書データ群から取得する。質問「フランスの首都は？」からの、キーワード「フランス」、「首都」の取り出しには、形態素解析技術などを使う。質問応答装置 1 が質問作成部 1 3 を省略し、解答候補抽出部 1 4 がキーワード増加部 1 2 によって出力された出力キーワードにより構成される質問に対する解答候補を抽出する構成を採る場合には、上記の質問「フランスの首都は？」からキーワード「フランス」、「首都」を形態素解析技術を用いて取り出す必要はなく、解答候補抽出部 1 4 は、キーワード増加部 1 2 によって出力された出力キーワード「フランス」、「首都」をそのまま用いて、それらを含む文書を上記文書データ群から取得する。

10

【 0 1 9 4 】

解答候補抽出部 1 4 は、キーワード「フランス」、「首都」を含む文書中の言語表現を、質問の答えの表現の候補として取り出す。この表現の取り出しには、例えば、前述の非特許文献 1 に記載された解の抽出の処理を用いる。

【 0 1 9 5 】

取り出された答えの表現の候補を、例えば、 $Score_{near}(c)$ の値の大きい順に並び替え、その値の上位何個かの候補を取り出し、その候補（候補 1, 候補 2, …）について、

20

問題『質問「フランスの首都は？」で答え「候補 1」』

問題『質問「フランスの首都は？」で答え「候補 2」』

・・・

を作成する。ここで、「候補 1」、「候補 2」は、上記の質問の「答え」の表現の候補を示している。

【 0 1 9 6 】

作成された問題（質問 - 答えの対）について、前述した機械学習の手法を適用し、「正解となりやすい」かの度合いが最も大きいときの、問題（質問 - 答えの対）における「答え」を、解答候補として、解答表出力部 1 5 に対して出力する。

【 0 1 9 7 】

30

解答表出力部 1 5 は、解答表において、質問「フランスの首都は？」に対する解答が格納される桁目（例えば、データ項目「フランス」に対応する行とデータ項目「首都」に対応する列とが交差する桁目）に、対応する解答候補を格納する。

【 0 1 9 8 】

（解答候補の抽出手法 2）

< 問題の構成 >

解答候補抽出部 1 4 は、

問題『質問「X 1 の Y 1 は？」』 - - - 解答「地名」

問題『質問「X 2 の Y 2 は？」』 - - - 解答「地名」

問題『質問「X 3 の Y 3 は？」』 - - - 解答「人名」

問題『質問「X 4 の Y 4 は？」』 - - - 解答「数値」

40

・・・

という、問題と解答の対を多数作成する。

素性としては、

X_i, Y_i の単語自体

X_i, Y_i の単語の意味クラス

などが考えられる。

問題構成と素性の定義をすれば、あとは機械学習の手法で扱える。

【 0 1 9 9 】

< 問題や素性の具体例 >

50

問題の具体例：

問題『質問「日本の首都は？」』 - - - 解答「地名」

問題『質問「日本の首相は？」』 - - - 解答「人名」

問題『質問「日本の面積は？」』 - - - 解答「数値」

素性の具体例：

問題『質問「日本の首都は？」』 - - - 解答「地名」の場合、

・ X i の単語自体：日本

・ Y i の単語自体：首都

・ X i の意味クラス：1 2 5 9 0（地名のクラス）

・ Y i の意味クラス：1 2 5 4 0（都市集落のクラス）

（意味クラスとして分類語彙表の最初の5桁を利用）

もっと多くの事例で学習すると、例えば、

Y i の単語自体：首都

だと、

解答「地名」

となるように学習し、

Y i の単語自体：首相

だと、

解答「人名」、

Y i の単語自体：面積

だと、

解答「数値」、

といったことを、解答候補抽出部14が学習し、その学習結果を学習DB17内に蓄積する。

【0200】

そして、解答候補抽出部14は、学習DB17内に蓄積された学習結果を用いて、解答を判断する。

例えば、新しい問題：

『質問「フランスの首都は？」』についての解答は、

Y i の単語自体：首都

なので、「地名」と判断する。

【0201】

<具体例>

例えば、質問作成部13によって作成された質問、または、質問応答装置1が質問作成部13を備えない構成を採るときはキーワード増加部12によって出力された出力キーワードによって構成される質問「フランスの首都は？」を解く場合について説明する。

【0202】

まず、解答候補抽出部14は、前述した機械学習の手法を利用して、

問題「フランスの首都は？」について、

解答が「地名」という結果を取得する。

【0203】

取得された解答「地名」を、解答表の桁目に格納する解答候補を抽出する際の解答タイプとして利用する。

【0204】

すなわち、解答候補抽出部14は、新聞記事データ・百科事典データなどの文書データ群から、質問作成部13が作成した質問「フランスの首都は？」を構成するキーワード（「フランス」、「首都」）を含む文書を取り出し、取り出された文書に含まれる言語表現のうち、上記解答タイプに適合するものを解答候補として解答表出力部15に対して出力する。解答候補抽出部14は、質問応答装置1が質問作成部13を備えない構成を採るときはキーワード増加部12によって出力された出力キーワード（「フランス」、「首都」

10

20

30

40

50

)を含む文書を上記文書データ群から取り出し、取り出された文書に含まれる言語表現のうち、上記解答タイプに適合するものを、出力キーワード「フランス」と「首都」とによって構成される質問「フランスの首都は？」に対する解答候補として解答表出力部15に対して出力する。

【0205】

解答表出力部15は、解答表において、質問「フランスの首都は？」に対する解答が格納される桁目(例えば、データ項目「フランス」に対応する行とデータ項目「首都」に対応する列とが交差する桁目)に、対応する解答候補を格納する。

【0206】

なお、本発明においては、例えば、「XのYは？」という質問に対する解答候補を抽出する際に、例えば、解答候補抽出部14が、機械学習の手法を用いるのではなく、新聞記事データ・百科事典データなどの大量の文書データ群(図示を省略)からキーワード「X」とキーワード「Y」を含む記事群を取り出し、その取り出した記事群の言語表現のうち、上記文書データ群中に出現する頻度が所定の閾値以上のものを解答候補として出力する構成を採ることもできる。また、本発明の実施の形態においては、上記取り出した記事群の言語表現について、上記文書データ群中に出現する頻度の高い順に所定の個数取り出して、解答候補として出力する構成を採ることもできる。

10

【0207】

ここで、上記の解答候補抽出部14による、解答タイプを用いた解答候補の出力の際には、非特許文献1の説明において述べた固有表現抽出技術を用いる。固有表現とは、人名、地名、組織名などの固有名詞、金額などの数値表現といった、特定の事物・数量を意味する言語表現のことで、固有表現抽出とは、そういった固有表現を文章中から計算機で自動で抽出する技術である。例えば、「日本の首相は小泉純一郎である」という文に対して固有表現抽出を行なうと、固有表現の「日本」と「小泉純一郎」が地名、人名として、抽出される。本発明の実施の形態においては、解答候補抽出部14が、抽出された固有表現が上記解答タイプに適合するかを判断し、適合する固有表現を、解答候補として出力する。

20

【0208】

以下に、固有表現抽出の一般的な手法の例について説明する。

(1) 機械学習を用いる手法

30

機械学習を用いて固有表現を抽出する手法がある(例えば、以下の文献(12)参照)

【0209】

文献(12): 浅原正幸, 松本裕治, 日本語固有表現抽出における冗長的な形態素解析の利用情報処理学会自然言語処理研究会 NL153-7 2002

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

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

40

50

す 0
。 0

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

【 0 2 1 0 】

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

10

【 0 2 1 1 】

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

日 B - LOCATION

の部分は、

日本 - B 名詞 - B

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

20

【 0 2 1 2 】

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

E O S

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

30

【 0 2 1 3 】

なお、例えば、上記の文献(12)では、素性として、入力文を構成する文字の、文字自体(例えば、「小」という文字)、字種(例えば、ひらがなやカタカナ等)、品詞情報、タグ情報(例えば、「B - PERSON」等)を利用している。

【 0 2 1 4 】

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

40

【 0 2 1 5 】

固有表現抽出には、上記の手法の他にも種々の手法がある。例えば、最大エントロピーモデルと書き換え規則を用いて固有表現を抽出する手法がある(文献(13)参照)。

【 0 2 1 6 】

文献(13): 内元清貴, 馬青, 村田真樹, 小作浩美, 内山将夫, 井佐原均, 最大エントロピーモデルと書き換え規則に基づく固有表現抽出, 言語処理学会誌, Vol.7, No.2, 2000 参照)。

【 0 2 1 7 】

また、例えば、以下の文献(14)に、サポートベクトルマシンを用いて日本語固有表現抽出を行う手法について記載されている。

【 0 2 1 8 】

50

文献(14)：山田寛康，工藤拓，松本裕治，Support Vector Machineを用いた日本語固有表現抽出，情報処理学会論文誌，Vol.43，No.1"，2002

(2) 形態素解析を用いる手法

形態素解析システム(例えば、前述したChaSen)を用いれば、入力された日本語を単語に分割することができる。

【0219】

例えば、ChaSenは、前述したように、日本語文を分割し、さらに、各単語の品詞も推定してくれる。例えば、「学校へ行く」を入力すると以下の結果を得ることができる。

【0220】

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

10

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

【0221】

具体的には、

入力：

日本の首都は東京です

出力：

日本ニッポン日本名詞 - 固有名詞 - 地域 - 国

のノの助詞 - 連体化

首都シュト首都名詞 - 一般

はハは助詞 - 係助詞

東京トウキョウ東京名詞 - 固有名詞 - 地域 - 一般

ですデスです助動詞特殊・デス基本形

EOS

は chasen の出力であり、名詞 - 固有名詞 - 地域という品詞が出力される。

このシステムを使って、例えば地名の固有表現を取り出すことができる。

【0222】

また、

入力：

村山首相が言った

出力：

村山ムラヤマ村山名詞 - 固有名詞 - 人名 - 姓

首相シュショウ首相名詞 - 一般

がガが助詞 - 格助詞 - 一般

言っイッ言う動詞 - 自立五段・ワ行促音便連用タ接続

たタた助動詞特殊・タ基本形

EOS

も chasen の出力であるが、これだと名詞 - 固有名詞 - 人名という品詞が出力される。このシステムを使って、例えば人名の固有表現を取り出すことができる。

(3) 作成したルールを用いる手法

人手でルールを作って固有表現を取り出すという方法もある。

【0223】

例えば、

名詞 + 「さん」だと人名とする

名詞 + 「首相」だと人名とする

名詞 + 「町」だと地名とする

名詞 + 「市」だと地名とする

50

などである。

【0224】

図5は、本発明の第1の実施の形態における質問応答処理フローの一例を示す図である。キーワード入力部11に、第1のキーワードと第2のキーワードを入力キーワードとして入力する(ステップS1)。例えば、第1のキーワード「日本」と、第2のキーワード「面積」とを入力する。

【0225】

キーワード増加部12のパターン抽出部121で、入力キーワードをキーワード抽出用DB16で全文検索し、入力キーワードの周辺に出現したパターンを c_i として抽出する(ステップS2)。周辺に出現するパターンの定義は適宜行なう。パターン c_i の抽出は、第1のキーワードと第2のキーワードそれぞれについて行う。

10

【0226】

キーワード増加部12のキーワード抽出部122で、パターン抽出部121で抽出したパターン c_i をキーワード抽出用DB16で全文検索し、パターン c_i によって抽出される表現 exp を抽出すると同時に、抽出した表現 exp をScoreの値の大きい順にソートし、キーワードとして出力する(ステップS3)。ステップS3の処理によって、例えば、第1のキーワードが、「日本」、「アメリカ」、「ドイツ」という3つの第3のキーワードに増加し、第2のキーワードが、「面積」、「人口」、「緯度」の3つの第4のキーワードに増加する。

【0227】

20

次に、質問作成部13が、出力されたキーワードにより構成される質問を作成する(ステップS4)。ステップS4においては、第3のキーワードと第4のキーワードとにより構成される質問を作成する。例えば、質問作成部13は、第3のキーワード「アメリカ」と第4のキーワード「人口」とにより構成される質問「アメリカの人口は？」を作成する。質問応答装置1が質問作成部13を備えない構成を採るときは、上記ステップS4の処理は、省略される。

【0228】

次に、解答候補抽出部14は、作成された各質問に対する解答候補を、上述した機械学習の手法を用いて抽出する(ステップS5)。質問応答装置1が質問作成部13を備えない構成を採るときは、上記ステップS5において、解答候補抽出部14は、キーワード増加部12によって出力されたキーワードによって構成される質問に対する解答候補を、機械学習の手法を用いて抽出する。そして、解答表出力部15が、解答表を出力する(ステップS6)。

30

【0229】

図6は、本発明の第1の実施の形態における質問応答装置の構成の別の例を示す図である。質問応答装置10は、入力されたキーワードを増加し、増加したキーワードにより構成される質問に対する解答を出力する装置である。図6中に示す質問応答装置10が備える構成要素のうち、図1に示す質問応答装置1が備える構成要素と同一の符号が付けられたものは、当該質問応答装置1が備える構成要素と同様の機能を有する。

【0230】

40

本発明の実施の形態においては、図6に示す構成から質問作成部13を省略し、解答候補抽出部14が、キーワード増加部60によって出力された第3のキーワードと第4のキーワードとによって構成される質問に対する解答候補を抽出し、出力する構成を採ってもよい。

【0231】

質問応答装置10のキーワード増加部60は、キーワード入力部11に入力されたキーワードを増加させる。すなわち、キーワード増加部60は、例えば、キーワード入力部11に入力された第1のキーワードに基づいて、第1のキーワードの数より多い第3のキーワードを出力する。また、キーワード増加部60は、例えば、キーワード入力部11に入力された第2のキーワードに基づいて、第2のキーワードの数より多い第4のキーワード

50

を出力する。

【0232】

単語データベース(DB)61には、単語と単語の分野との対応情報が格納されている。例えば、図7に示すような、単語と単語の分野との対応情報が格納されている。例えば、「国名」という分野に対応する単語として、日本、アメリカ、ドイツ、・・・といった単語が格納されている。

【0233】

また、シソーラスデータベース(DB)62には、意味的類似による単語の分類情報であるシソーラスデータが格納されている。例えば、シソーラスDB62には、図8に示すような、単語と単語に振られた10桁の数字(分類番号)との対応情報がシソーラスデータとして格納されている。図8に示す例では、シソーラスデータが分類語彙表の形式で示されている。

10

【0234】

なお、分類語彙表とは、一般に、単語を意味に基づいて整理した表であり、各単語に対して分類番号という数字が付与されている。この10桁の分類番号は、7レベルの階層構造を示しており、上位5レベルは分類番号の最初の5桁で表現され、6レベル目は次の2桁、最下層のレベルは最後の3桁で表現されている。

【0235】

類似度算出部100は、シソーラスDB62中のシソーラスデータに基づいて、キーワード入力部11に入力されたキーワードとシソーラスデータ中の単語との類似度を算出する。キーワード抽出部101は、例えば、算出された類似度が予め定めた閾値以上の単語をキーワードとして抽出し、出力する。また、キーワード抽出部101は、例えば、算出された類似度が大きい順に所定の個数の単語をシソーラスデータ中から取り出して、キーワードとして出力する構成を採ることもできる。

20

【0236】

本発明の実施の形態においては、キーワード抽出部101は、単語データDB61中に格納された、単語と単語の分野との対応情報に基づいて、キーワード入力部11に入力されたキーワードと同じ分野の単語をキーワードとして抽出し、出力する構成を採ることもできる。

【0237】

上記の質問応答装置10を用いた場合の質問応答処理フローは、図5に示す質問応答処理フローと、ステップS2、ステップS3の処理が異なる以外は、同様である。質問応答装置10を用いた場合の質問応答処理フローの一例においては、図5のステップS2およびステップS3の代わりに、キーワード増加部60のキーワード抽出部101で、キーワード入力部11に入力されたキーワードと同じ分野の単語を単語データDB61中から抽出し、キーワードとして出力する。

30

【0238】

例えば、キーワード入力部11に第1のキーワード「日本」が入力されたとすると、キーワード抽出部101は、図7に示す単語データDB61から、単語「日本」が対応する「国名」という分野に属する(対応する)単語である「日本」、「アメリカ」、「ドイツ」、・・・を抽出し、第3のキーワードとして出力する。また、例えば、キーワード入力部11に第2のキーワード「面積」が入力されたとすると、キーワード抽出部101は、図7に示す単語データDB61から、単語「面積」が対応する「数値表現」という分野に属する(対応する)単語である「面積」、「人口」、「緯度」、・・・を抽出し、第4のキーワードとして出力する。

40

【0239】

また、質問応答装置10を用いた場合の質問応答処理フローの別の例においては、図5のステップS2およびステップS3の代わりに、例えば、キーワード増加部60の類似度算出部100が、キーワード入力部11に入力されたキーワードとシソーラスDB62中の単語との類似度を算出し、キーワード増加部60のキーワード抽出部101が、算出さ

50

れた類似度が予め定めた閾値以上の単語をキーワードとして出力する。

【0240】

なお、例えば、キーワード抽出部101は、算出された類似度が大きい順に所定の個数の単語をシソーラスデータ中から取り出して、キーワードとして出力する構成を採ることもできる。

【0241】

類似度算出部100は、入力されたキーワードとシソーラスDB62中の単語との類似度を、例えば以下のようにして算出する。図8に示すシソーラスDB62内に格納されたシソーラスデータ(分類語彙表)中の各単語に振られた、10桁の分類番号における各桁の数字の一致の割合を用いて、類似度を求める。すなわち、例えば、分類語彙表中の各単語に振られた分類番号について、キーワード入力部11に入力されたキーワードと同一の単語に振られた分類番号との間での、各桁の数字の一致の割合を算出し、算出された値を類似度とする。なお、例えば、分類番号の6桁目と7桁目、および、8桁目と9桁目と10桁目は、それぞれ連続した1つの数字として考える。

10

【0242】

例えば、キーワード入力部11に第1のキーワードとして入力されたキーワードが「日本」である場合、図8に示す分類語彙表中の単語「日本」と「アメリカ」には、それぞれ以下のような分類番号が振られている。以下では、分類番号の上位5レベルと、6レベル目と、最下層のレベルとの間を空白で区切って示す。

【0243】

日本 : 1 2 5 9 0 0 1 0 1 2
アメリカ : 1 2 5 9 0 0 4 1 9 2

例えば、両単語の分類番号の上位5レベルにおいて、最初の5桁が一致するので、算出されるキーワード「日本」と分類語彙表中の単語「アメリカ」との類似度は、類似度5である。

20

【0244】

また、例えば、キーワード入力部11に第2のキーワードとして入力されたキーワードが「面積」である場合、分類語彙表中の単語「面積」と「人口」には、それぞれ以下のような分類番号が振られている。

【0245】

面積 : 1 2 6 3 0 1 3 0 1 5
人口 : 1 2 6 3 0 1 0 0 1 2

例えば、両単語の分類番号の上位5レベルにおいて、最初の5桁が一致するので、算出されるキーワード「面積」と分類語彙表中の単語「人口」との類似度は、類似度5である。

30

【0246】

また、例えば、キーワード入力部11に第2のキーワードとして入力されたキーワードが「人口」である場合、分類語彙表中の単語「人口」と「緯度」には、それぞれ以下のような分類番号が振られている。

【0247】

人口 : 1 2 6 3 0 1 0 0 1 2
緯度 : 1 2 6 3 0 1 0 0 1 5

例えば、両単語の分類番号の上位5レベルにおいて、最初の5桁が一致し、また、6レベル目の2桁の数字「10」が一致するので、算出されるキーワード「人口」と分類語彙表中の単語「緯度」との類似度は、類似度7である。

40

【0248】

また、例えば、キーワード入力部11に第2のキーワードとして入力されたキーワードが「人口」である場合、分類語彙表中の単語「人口」と「アメリカ」には、それぞれ以下のような分類番号が振られている。

【0249】

50

人口 : 1 2 6 3 0 1 0 0 1 2

アメリカ : 1 2 5 9 0 0 4 1 9 2

例えば、両単語の分類番号の上位5レベルにおいて、最初の2桁が一致するため、算出されるキーワード「人口」と分類語彙表中の単語「アメリカ」との類似度は、類似度2である。

【0250】

図9は、本発明の第2の実施の形態における質問応答装置の構成の一例を示す図である。第2の実施の形態においては、例えば、第1のキーワード「日本」と、第2のキーワード「首都」+疑問代名詞「はどこですか?」が入力されると、第1のキーワード「日本」に基づいて、第1のキーワードを、例えば「日本」、「アメリカ」、「ドイツ」の3つに増加させる。そして、増加後の第1のキーワードと、第2のキーワード「首都」+疑問代名詞「はどこですか?」により構成される、例えば「日本の首都はどこですか?」、「アメリカの首都はどこですか?」、「ドイツの首都はどこですか?」という各質問に対する解答を出力する。

10

【0251】

質問応答装置2は、入力されたキーワードを増加し、増加したキーワードにより構成される質問に対する解答を出力する装置である。図9に示す質問応答装置2の構成要素のうち、キーワード入力部11、解答表出力部15、キーワード抽出用DB16、パターン抽出部121、キーワード抽出部122は、それぞれ、図1に示す質問応答装置1の、同符号の構成要素と同様である。本発明の実施の形態においては、図9に示す構成から後述する質問作成部23を省略し、解答候補抽出部24が、キーワード増加部18によって出力されたキーワードによって構成される質問に対する解答候補を抽出し、出力する構成を採ってもよい。

20

【0252】

キーワード入力部11には、キーワードが入力される。例えば、第1のキーワード「日本」と第2のキーワード「首都」が入力される。疑問代名詞入力部21には、キーワード入力部11に入力された第2のキーワードに対応付けられた疑問代名詞が入力される。例えば、「はどこですか?」という疑問代名詞が入力される。この他、疑問代名詞入力部21に入力される疑問代名詞として、例えば、「は何時ですか?」、「は誰ですか?」などが挙げられる。なお、疑問代名詞入力部21に入力される疑問代名詞は、ユーザの指定入力に基づいて入力されるものであってもよいし、また、質問応答装置2とは別のコンピュータによって入力されるものであってもよい。

30

【0253】

解答タイプ推定部22は、疑問代名詞入力部21に入力された疑問代名詞に基づいて、後述する質問作成部23によって作成される質問、または、質問応答装置2が質問作成部23を備えない構成を採るときは、後述するキーワード増加部18によって出力されたキーワードによって構成される質問に対する解答候補の言語表現のタイプである解答タイプを推定する。例えば、入力された疑問代名詞が「はどこですか?」である場合には、解答タイプは「固有名詞(地名)」であると推定する。本発明の実施の形態においては、解答タイプ推定部22は、疑問代名詞入力部21に入力された疑問代名詞ではなく、予め定められた疑問代名詞に基づいて、上記解答タイプを推定してもよい。

40

【0254】

キーワード増加部18は、キーワード抽出技術を用いて、入力された第1のキーワードと同じ分野のキーワードをキーワード抽出用DB16から抽出して、第1のキーワードを増加させ、第3のキーワードとして出力する。第2の実施の形態では、キーワード増加部18は、第2のキーワード(例えば、「首都」)については増加させずに、質問作成部23に対して出力する。質問応答装置2が質問作成部23を備えない構成を採るときは、キーワード増加部18は、第3のキーワードと第2のキーワードを解答候補抽出部24に対して出力する。

【0255】

50

質問作成部 23 は、キーワード増加部 18 の処理によって出力された第 3 のキーワードと、第 2 のキーワードと、疑問代名詞入力部 21 に入力された疑問代名詞（または予め定められた疑問代名詞）とに基づいて、複数の質問を作成する。

【0256】

知識データベース (DB) 25 には、解答候補の検索対象となる文書データ群が蓄積される。蓄積される文書データ群としては、例えば、新聞記事データ・百科事典データなどの文書データ群が挙げられる。

【0257】

解答候補抽出部 24 は、知識 DB 25 から、質問作成部 23 によって作成された各質問を構成するキーワード（または、キーワード増加部 18 から出力された第 3 のキーワードと第 2 のキーワード）を含む文書データを検索し、この検索処理で抽出された文書データから、解答タイプ推定部 22 によって推定された解答タイプに適合する言語表現を、解答候補として抽出する。

10

【0258】

解答表出力部 15 は、抽出された各解答候補が質問と対応付けられた表を解答表として出力する。例えば、図 10 に示すような解答表を出力する。

【0259】

図 10 に示す解答表においては、例えば、「日本の首都はどこですか?」という質問に対する解答として、データ項目「日本」に対応する行とデータ項目「首都」に対応する列とが交差する桁目に、「東京」が格納され、「アメリカの首都はどこですか?」という質問に対する解答として、データ項目「アメリカ」に対応する行とデータ項目「首都」に対応する列とが交差する桁目に、「ワシントン」が格納され、「ドイツの首都はどこですか?」という質問に対する解答として、データ項目「アメリカ」に対応する行とデータ項目「首都」に対応する列とが交差する桁目に、「ベルリン」が格納される。

20

【0260】

図 11 は、本発明の第 2 の実施の形態における質問応答処理フローの一例を示す図である。キーワード入力部 11 に、第 1 のキーワードと第 2 のキーワードを入力キーワードとして入力する (ステップ S11)。例えば、第 1 のキーワード「日本」と第 2 のキーワード「首都」が入力される。また、疑問代名詞入力部 21 に、第 2 のキーワードに対応付けられた疑問代名詞が入力される (ステップ S12)。例えば、第 2 のキーワード「首都」に対応付けられた疑問代名詞「はどこですか?」が入力される。

30

【0261】

キーワード増加部 18 のパターン抽出部 121 で、第 1 のキーワードをキーワード抽出用 DB 16 で全文検索し、第 1 のキーワードの周辺に出現したパターンを c_i として抽出する (ステップ S13)。周辺に出現するパターンの定義は適宜行なう。

【0262】

キーワード増加部 18 のキーワード抽出部 122 で、パターン抽出部 121 で抽出したパターン c_i をキーワード抽出用 DB 16 で全文検索し、パターン c_i によって抽出される表現 $e x p$ を抽出すると同時に、抽出した表現 $e x p$ を $S c o r e$ の値の大きい順にソートし、第 3 のキーワードとして出力する (ステップ S14)。ステップ S14 の処理によって、例えば、第 1 のキーワードが、「日本」、「アメリカ」、「ドイツ」という 3 つの第 3 のキーワードに増加する。

40

【0263】

解答タイプ推定部 22 が、疑問代名詞入力部 21 に入力された疑問代名詞に基づいて、解答タイプを推定する (ステップ S15)。例えば、入力された疑問代名詞が「はどこですか?」である場合には、解答タイプ推定部 22 は、解答タイプが「固有名詞 (地名)」であると推定する。

【0264】

質問作成部 23 が、疑問代名詞入力部 21 に入力された疑問代名詞を用いて、第 3 のキーワードと第 2 のキーワードとにより構成される質問を作成する (ステップ S16)。例

50

えば、質問作成部 23 は、第 3 のキーワード「アメリカ」と第 2 のキーワード「首都」とにより構成される質問「アメリカの首都はどこですか？」を作成する。質問応答装置 2 が質問作成部 23 を備えない構成を採るときは、上記ステップ S 16 の処理は、省略される。

【0265】

次に、解答候補抽出部 24 は、作成された各質問に対する解答候補を抽出する（ステップ S 17）。すなわち、解答候補抽出部 24 は、知識 DB 25 から、質問作成部 23 によって作成された各質問を構成するキーワードを含む文書データを検索し、この検索処理で抽出された文書データから、解答タイプ推定部 22 によって推定された解答タイプに適合する言語表現を、解答候補として抽出する。質問応答装置 2 が質問作成部 23 を備えない構成を採るときは、上記ステップ S 17 において、解答候補抽出部 24 は、知識 DB 25 から、キーワード増加部 18 によって出力されたキーワードを含む文書データを検索し、この検索処理で抽出された文書データから、解答タイプ推定部 22 によって推定された解答タイプに適合する言語表現を、キーワード増加部 18 によって出力されたキーワードによって構成される質問に対する解答候補として抽出する。そして、解答表出力部 15 が、解答表を出力する（ステップ S 18）。例えば、上述した図 10 に示すような解答表が出力される。

10

【0266】

図 12 は、本発明の第 2 の実施の形態の変形例 1 の構成例を示す図である。質問応答装置 20 は、入力されたキーワードを増加し、増加したキーワードにより構成される質問に対する解答を出力する装置である。図 12 中に示す質問応答装置 10 が備える構成要素のうち、図 6 に示す質問応答装置 10 が備える構成要素または図 9 に示す質問応答装置 2 が備える構成要素と同一の符号が付けられたものは、当該質問応答装置 10 または質問応答装置 2 が備える構成要素と同様の機能を有する。

20

【0267】

質問応答装置 20 のキーワード増加部 63 は、キーワード入力部 11 に入力された第 1 のキーワードを増加させて、第 3 のキーワードとして出力する。また、キーワード入力部 11 に入力された第 2 のキーワードについては、増加させずに、質問作成部 23 に対して出力する。本発明の実施の形態においては、図 12 に示す構成から質問作成部 23 を省略し、解答候補抽出部 24 が、キーワード増加部 63 によって出力された第 3 のキーワードと第 2 のキーワードとによって構成される質問に対する解答候補を抽出し、出力する構成を採ってもよい。

30

【0268】

上記の質問応答装置 20 を用いた場合の質問応答処理フローは、図 11 に示す質問応答処理フローと、ステップ S 13、ステップ S 14 の処理が異なる以外は、同様である。質問応答装置 10 を用いた場合の質問応答処理フローの一例においては、図 11 のステップ S 13 およびステップ S 14 の代わりに、キーワード増加部 63 のキーワード抽出部 101 で、キーワード入力部 11 に入力された第 1 のキーワードと同じ分野の単語を単語データ DB 61 中から抽出し、第 3 のキーワードとして出力する。

40

【0269】

また、質問応答装置 10 を用いた場合の質問応答処理フローの別の例においては、図 11 のステップ S 13 およびステップ S 14 の代わりに、キーワード増加部 63 の類似度算出部 100 が、キーワード入力部 11 に入力された第 1 のキーワードとシソーラス DB 62 中の単語との類似度を算出し、キーワード増加部 63 のキーワード抽出部 101 が、算出された類似度が予め定めた閾値以上の単語を第 3 のキーワードとして出力する。

【0270】

また、キーワード抽出部 101 は、例えば、算出された類似度が大きい順に所定の個数の単語をシソーラスデータ中から取り出して、第 3 のキーワードとして出力する構成を採ることもできる。

【0271】

50

本発明の第2の実施の形態の変形例2においては、図9に示す質問応答装置2または図12に示す質問応答装置20において、疑問代名詞入力部21には、キーワード入力部11に入力されるキーワードと対応付けられていない疑問代名詞が入力される。質問応答装置2のキーワード増加部18（または質問応答装置20のキーワード増加部63）は、キーワード入力部11に入力された第1のキーワードに基づいて、第3のキーワードを出力キーワードとして出力し、キーワード入力部11に入力された第2のキーワードに基づいて、第4のキーワードを出力キーワードとして出力する。解答タイプ推定部22は、疑問代名詞入力部21に入力された疑問代名詞に基づいて、解答タイプを推定する。そして、解答候補抽出部24は、知識DB25から、キーワード増加部18によって出力されたキーワードを含む文書データを検索し、この検索処理で抽出された文書データから、解答タイプ推定部22によって推定された解答タイプに適合する言語表現を、キーワード増加部18によって出力されたキーワードによって構成される質問に対する解答候補として抽出する。そして、解答表出力部15が、解答表を出力する。なお、上記の本発明の第2の実施の形態の変形例2においては、解答タイプ推定部22は、疑問代名詞入力部21に入力された疑問代名詞ではなく、予め定められた疑問代名詞に基づいて解答タイプを推定する構成を採ってもよい。

10

【0272】

図13は、本発明の第3の実施の形態における質問応答装置の構成の一例を示す図である。第3の実施の形態では、第1の実施の形態のような機械学習の手法を用いるのではなく、入力された解答タイプ（または予め定められた解答タイプ）を用いて解答候補を抽出する。

20

【0273】

第3の実施の形態においては、例えば、第1のキーワード「日本」と第2のキーワード「首都」と、解答タイプ「固有名詞（地名）」が入力されると、第1のキーワード「日本」に基づいて、第1のキーワードを、例えば「日本」、「アメリカ」、「ドイツ」という3つの第3のキーワードに増加させる。また、第2のキーワード「首都」に基づいて、第2のキーワードを、例えば「首都」、「旧首都」、「最南端都市」という3つの第4のキーワードに増加させる。

【0274】

そして、増加後の第3のキーワードと第4のキーワードとの組み合わせにより構成される、例えば「日本の首都は?」、「アメリカの旧首都は?」、「ドイツの最南端都市は?」・・・といった各質問に対する解答を出力する。より具体的には、後述するように、「日本の首都は?」という質問を構成する第3のキーワード「日本」と第4のキーワード「首都」を、解答候補の検索対象となる文書データ群から検索し、両キーワードを含む文書中の言語表現を解答候補として抽出するとともに、抽出された解答候補のうち解答タイプ「固有名詞（地名）」に適合するものを解答として出力する。

30

【0275】

質問応答装置3は、入力されたキーワードを増加し、増加したキーワードにより構成される質問に対する解答を出力する装置である。図13に示す質問応答装置3の構成要素のうち、キーワード入力部11、キーワード増加部12、質問作成部13、解答表出力部15、キーワード抽出用DB16、パターン抽出部121、キーワード抽出部122は、それぞれ、図1に示す質問応答装置1の、同符号の構成要素と同様であり、解答候補抽出部24、知識DB25は、図9に示す質問応答装置2の、同符号の構成要素と同様である。本発明の実施の形態においては、図13に示す構成から質問作成部13を省略し、解答候補抽出部24が、キーワード増加部12によって出力されたキーワードによって構成される質問に対する解答候補を抽出し、出力する構成を採ってもよい。

40

【0276】

キーワード入力部11には、キーワードが入力される。例えば、第1のキーワード「日本」と第2のキーワード「首都」が入力される。解答タイプ入力部31には、質問作成部13によって作成される質問、または、質問応答装置3が質問作成部13を省略する構成

50

を採るときは、キーワード増加部 1 2 によって出力されるキーワードによって構成される質問に対する解答候補の解答タイプが入力される。例えば、「固有名詞(地名)」という解答タイプが入力される。

【0277】

この他、解答タイプ入力部 3 1 に入力される解答タイプとして、例えば、「固有名詞(数値)」、「固有名詞(人名)」、「カタカナ表現」(カタカナだけで表現されるもの)、「名詞」、「動詞」などが挙げられる。なお、解答タイプ入力部 3 1 に入力される解答タイプは、ユーザの指定入力に基づいて入力されるものであってもよいし、また、質問応答装置 3 とは別のコンピュータによって入力されるものであってもよい。

【0278】

キーワード増加部 1 2 は、図 1 を参照して説明したように、キーワード抽出技術を用いて、入力された各キーワードと同じ分野のキーワードをキーワード抽出用 DB 1 6 から抽出して、キーワードを増加させる。

【0279】

例えば、キーワード増加部 1 2 は、入力された第 1 のキーワードに基づいて、第 1 のキーワードの数より多い第 3 のキーワードを出力する。また、例えば、キーワード増加部 1 2 は、入力された第 2 のキーワードに基づいて、第 2 のキーワードの数より多い第 4 のキーワードを出力する。

【0280】

質問作成部 1 3 は、第 3 のキーワードと第 4 のキーワードとによって構成される質問を複数作成する。解答候補抽出部 2 4 は、知識 DB 2 5 から、質問作成部 1 3 によって作成された各質問を構成するキーワード(または、キーワード増加部 1 2 によって出力された第 3 のキーワードと第 4 のキーワード)を含む文書を検索し、この検索処理で抽出された文書に含まれる言語表現のうち、解答タイプ入力部 3 1 に入力された解答タイプ(または予め定められた解答タイプ)に適合する言語表現を、解答候補として抽出する。

【0281】

解答表出力部 1 5 は、抽出された各解答候補が質問と対応付けられた表を解答表として出力する。

【0282】

図 1 4 は、本発明の第 3 の実施の形態における質問応答処理フローの一例を示す図である。キーワード入力部 1 1 に、第 1 のキーワードと第 2 のキーワードを入力キーワードとして入力する(ステップ S 2 1)。例えば、第 1 のキーワード「日本」と第 2 のキーワード「首都」が入力される。また、解答タイプ入力部 3 1 に、質問作成部 1 3 により作成される質問に対する解答候補の解答タイプを入力する(ステップ S 2 2)。例えば、解答タイプとして、「固有名詞(地名)」が入力される。なお、質問応答装置 3 が質問作成部 1 3 を備えない構成を採るときは、解答タイプ入力部 3 1 には、キーワード増加部 1 2 によって出力されるキーワードによって構成される質問に対する解答候補の解答タイプが入力される。

【0283】

キーワード増加部 1 2 のパターン抽出部 1 2 1 で、入力キーワードをキーワード抽出用 DB 1 6 で全文検索し、複数の入力キーワードの周辺に出現したパターンを c_i として抽出する(ステップ S 2 3)。周辺に出現するパターンの定義は適宜行なう。パターン c_i の抽出は、第 1 のキーワードと第 2 のキーワードそれぞれについて行う。

【0284】

キーワード増加部 1 2 のキーワード抽出部 1 2 2 で、パターン抽出部 1 2 1 で抽出したパターン c_i をキーワード抽出用 DB 1 6 で全文検索し、パターン c_i によって抽出される表現 $e x p$ を抽出すると同時に、抽出した表現 $e x p$ を $S c o r e$ の値の大きい順にソートし、キーワードとして出力する(ステップ S 2 4)。ステップ S 2 4 の処理によって、例えば、第 1 のキーワードが、「日本」、「アメリカ」、「ドイツ」という 3 つの第 3 のキーワードに増加する。また、第 2 のキーワードが、「首都」、「旧首都」、「最南端

10

20

30

40

50

都市」という3つの第4のキーワードに増加する。

【0285】

質問作成部23が、出力されたキーワードにより構成される質問を作成する(ステップS25)。ステップS25においては、出力された第3のキーワードと第4のキーワードとにより構成される質問を作成する。例えば、質問作成部23は、第3のキーワード「アメリカ」と第4のキーワード「首都」とにより構成される質問「アメリカの首都は?」を作成する。質問応答装置3が質問作成部13を備えない構成を採るときは、上記ステップS25の処理は、省略される。

【0286】

次に、解答候補抽出部24は、作成された各質問に対する解答候補を抽出する(ステップS26)。すなわち、解答候補抽出部24は、知識DB25から、質問作成部13によって作成された各質問を構成するキーワードを含む文書データを検索し、この検索処理で抽出された文書データから、解答タイプ入力部31に入力された解答タイプに適合する言語表現を、解答候補として抽出する。質問応答装置3が質問作成部13を備えない構成を採るときは、上記ステップS26において、解答候補抽出部24は、知識DB25から、キーワード増加部12によって出力されたキーワード(第3のキーワードと第4のキーワード)を含む文書データを検索し、この検索処理で抽出された文書データから、解答タイプ入力部31に入力された解答タイプに適合する言語表現を、キーワード増加部12によって出力された第3のキーワードと第4のキーワードとによって構成される質問に対する解答候補として抽出する。そして、解答表出力部15が、解答表を出力する(ステップS27)。例えば、図15に示すような解答表が出力される。

【0287】

図16は、本発明の第3の実施の形態の変形例1の構成例を示す図である。質問応答装置30は、入力されたキーワードを増加し、増加したキーワードにより構成される質問に対する解答を出力する装置である。図16中に示す質問応答装置30が備える構成要素のうち、図1に示す質問応答装置1が備える構成要素または図6に示す質問応答装置10または図13に示す質問応答装置3が備える構成要素と同一の符号が付けられたものは、当該質問応答装置1または質問応答装置10または質問応答装置3が備える構成要素と同様の機能を有する。

【0288】

本発明の実施の形態においては、図16に示す構成から質問作成部13を省略し、解答候補抽出部24が、キーワード増加部60によって出力された第3のキーワードと第4のキーワードとによって構成される質問に対する解答候補を抽出し、出力する構成を採ってもよい。

【0289】

上記の質問応答装置30を用いた場合の質問応答処理フローは、図14に示す質問応答処理フローと、ステップS23、ステップS24の処理が異なる以外は、同様である。質問応答装置30を用いた場合の質問応答処理フローの一例においては、図14のステップS23およびステップS24の代わりに、キーワード増加部60のキーワード抽出部101で、キーワード入力部11に入力された第1のキーワードと同じ分野の単語を単語データDB61中から抽出し、第3のキーワードとして出力する。また、キーワード抽出部101で、キーワード入力部11に入力された第2のキーワードと同じ分野の単語を単語データDB61中から抽出し、第4のキーワードとして出力する。

【0290】

また、質問応答装置30を用いた場合の質問応答処理フローの別の例においては、図14のステップS23およびステップS24の代わりに、キーワード増加部60の類似度算出部100が、キーワード入力部11に入力された第1のキーワードとシソーラスDB62中の単語との類似度を算出し、キーワード増加部60のキーワード抽出部101が、算出された類似度が予め定めた閾値以上の単語を第3のキーワードとして出力する。また、類似度算出部100が、キーワード入力部11に入力された第2のキーワードとシソーラ

10

20

30

40

50

スDB62中の単語との類似度を算出し、キーワード抽出部101が、算出された類似度が予め定めた閾値以上の単語を第4のキーワードとして出力する。

【0291】

また、キーワード抽出部101は、例えば、算出された類似度が大きい順に所定の個数の単語をシソーラスデータ中から取り出して、上記の第3のキーワード、第4のキーワードとして出力する構成を採ることもできる。

【0292】

本発明の第3の実施の形態の変形例2では、例えば、図13に示す質問応答装置3または図16に示す質問応答装置30において、解答タイプ入力部31には、キーワード増加部12（またはキーワード増加部60）によって出力される出力キーワードによって構成される質問に対する解答の候補の言語表現の類型である解答タイプであって、キーワード入力部11に入力された第2のキーワードに対応付けられた解答タイプが入力される。キーワード増加部12（またはキーワード増加部60）は、キーワード入力部11に入力された第1のキーワードに基づいて、第3のキーワードを出力キーワードとして出力し、キーワード入力部11に入力された第2のキーワードを出力キーワードとして出力する。解答候補抽出部24は、知識DB25から、キーワード増加部12（またはキーワード増加部60）によって出力されたキーワード（第3のキーワードと第2のキーワード）を含む文書データを検索し、この検索処理で抽出された文書データから、当該第2のキーワードに対応する解答タイプ入力部31に入力された解答タイプに適合する言語表現を、キーワード増加部12によって出力された第3のキーワードと当該第2のキーワードとによって構成される質問に対する解答候補として抽出する。そして、解答表出力部15が、解答表を出力する。

【0293】

なお、本発明の第3の実施の形態の変形例3では、解答候補抽出部24は、知識DB25から、キーワード増加部12（またはキーワード増加部60）によって出力されたキーワード（第3のキーワードと第2のキーワード）を含む文書データを検索し、この検索処理で抽出された文書データから、予め定められた、当該第2のキーワードに対応付けられた解答タイプに適合する言語表現を、キーワード増加部12によって出力された第3のキーワードと当該第2のキーワードとによって構成される質問に対する解答候補として抽出する構成を採ってもよい。

【0294】

本発明の第3の実施の形態の変形例3では、例えば、図13に示す質問応答装置3または図16に示す質問応答装置30において、キーワード入力部11には、第1のキーワードと、複数のグループによってグループ化された第2のキーワードとが入力される。例えば、第1のキーワード「日本」と、人名のグループに属する第2のキーワード「首相」、「市長」と、地名のグループに属する第2のキーワード「首都」、「旧首都」が入力される。解答タイプ入力部31には、キーワード入力部11に入力される第2のキーワードが属する各グループに対応付けられた解答タイプが入力される。例えば、人名のグループに対応する解答タイプとして、解答タイプ「人名」が入力され、地名のグループに対応する解答タイプとして、解答タイプ「地名」が入力される。

【0295】

キーワード増加部12（またはキーワード増加部60）は、第1のキーワードに基づいて、第3のキーワードを出力キーワードとして出力する。また、第2のキーワードに基づいて、当該第2のキーワードが属するグループ毎に、第4のキーワードを出力キーワードとして出力する。例えば、第3のキーワードとして、「日本」、「ドイツ」、「アメリカ」が出力される。また、例えば、人名のグループに属する第4のキーワードとして、「首相」、「市長」、「ノーベル賞受賞者」が出力され、地名のグループに属する第4のキーワードとして、「首都」、「旧首都」、「最南端都市」が出力される。

【0296】

解答候補抽出部24は、キーワード増加部12（またはキーワード増加部60）によ

10

20

30

40

50

て出力された第3のキーワードと、人名のグループに属する第4のキーワードとによって構成される質問に対する解答候補を、解答タイプ入力部31に入力された解答タイプ「人名」を用いて抽出する。例えば、「ドイツのノーベル賞受賞者は？」という質問に対する解答候補は、解答タイプ「人名」を用いて抽出される。また、解答候補抽出部24は、キーワード増加部12（またはキーワード増加部60）によって出力された第3のキーワードと、地名のグループに属する第4のキーワードとによって構成される質問に対する解答候補を、解答タイプ入力部31に入力された解答タイプ「地名」を用いて抽出する。例えば、「アメリカの首都は？」という質問に対する解答候補は、解答タイプ「地名」を用いて抽出される。そして、解答表出力部15が解答表を出力する。

【0297】

図17は、本発明の第4の実施の形態における質問応答装置の構成の一例を示す図である。第4の実施の形態では、キーワードの類似関係を用いて解答候補を抽出する。

【0298】

第4の実施の形態においては、例えば、第1のキーワード「日本」、「アメリカ」、・・・と第2のキーワード「面積」、「首都」、・・・と解答タイプ「固有名詞（数値）」、「固有名詞（地名）」、・・・が入力される。入力される解答タイプは、入力された第2のキーワードのそれぞれに対応付けられている。例えば、第2のキーワード「面積」に対応付けられた解答タイプは「固有名詞（数値）」であり、第2のキーワード「首都」に対応付けられた解答タイプは「固有名詞（地名）」である。

【0299】

第4の実施の形態では、例えば、入力された第1のキーワード「日本」、「アメリカ」、・・・に基づいて、第1のキーワードを多数の第3のキーワード（例えば「日本」、「アメリカ」、「ドイツ」、・・・）に増加させる。また、キーワード増加部12が第2のキーワード「面積」、「首都」、・・・に基づいて、第2のキーワードを多数の第4のキーワード（例えば「面積」、「首都」、「旧首都」、・・・）に増加させる。

【0300】

次に、第2のキーワード（と同一の第4のキーワード）のうち、第4のキーワードに類似するキーワードを、類似キーワードとして決定する。例えば、第4のキーワード「旧首都」に類似する第2のキーワード（と同一の第4のキーワード）「首都」を類似キーワードとして決定する。

【0301】

そして、第3のキーワードと第4のキーワードとの組み合わせにより構成される質問に対する解答の候補を、上記質問を構成する第4のキーワードに類似する類似キーワードに対応付けられている解答タイプを用いて抽出し、解答表を出力する。例えば、「日本の旧首都は？」という質問に対する解答の候補を、類似キーワード「首都」に対応付けられている解答タイプ「固有名詞（地名）」を用いて抽出し、解答表を出力する。

【0302】

質問応答装置4は、入力されたキーワードを増加し、増加したキーワードにより構成される質問に対する解答を出力する装置である。図17に示す質問応答装置4の構成要素のうち、キーワード入力部11、キーワード増加部12、解答表出力部15、キーワード抽出用DB16、パターン抽出部121、キーワード抽出部122は、それぞれ、図1に示す質問応答装置1の、同符号の構成要素と同様であり、知識DB25は、図9に示す質問応答装置2が備える知識DB25と同様であり、解答タイプ入力部31は、図13に示す質問応答装置3が備える解答タイプ入力部31と同様である。

【0303】

キーワード入力部11には、キーワードが入力される。例えば、第1のキーワード「日本」、「アメリカ」、・・・と第2のキーワード「面積」、「首都」、・・・が入力される。解答タイプ入力部31には、質問作成部42によって作成される質問、または、質問応答装置4が質問作成部42を省略する構成を採るときは、キーワード増加部12によって出力されるキーワードによって構成される質問に対する解答候補の解答タイプが入力され

10

20

30

40

50

る。入力される解答タイプは、特に、第2のキーワードに対応付けられている。

【0304】

例えば、解答タイプ入力部31には、第2のキーワード「面積」に対応して、「固有名詞(数値)」という解答タイプが入力され、第2のキーワード「首都」に対応して、「固有名詞(地名)」という解答タイプが入力される。

【0305】

キーワード増加部12は、図1を参照して説明したように、キーワード抽出技術を用いて、入力された各キーワードと同じ分野のキーワードをキーワード抽出用DB16から抽出して、キーワードを増加させる。キーワード増加部12の処理により、第1のキーワードから第3のキーワードが出力され、第2のキーワードから第4のキーワードが出力される。

10

【0306】

類似キーワード決定部41は、各第4のキーワードに類似する、キーワード入力部11に入力された第2のキーワード(と同一の第4のキーワード)を、類似キーワードとして決定する。類似キーワードの決定手法について以下に説明する。

【0307】

(共起ベクトルを用いる手法(1))

第4のキーワード毎に、キーワード増加部12が抽出したパターン c_i と共起してキーワード抽出用DB16中出现した回数を算出し、算出した回数を要素とするベクトル(以下、「共起ベクトル」という)を求める。

20

【0308】

例えば、キーワード増加部12におけるキーワード抽出処理において、第4のキーワード(1)がパターン c_1 と共起して出現した回数が0、パターン c_2 と共起して出現した回数が1、・・・、パターン c_n と共起して出現した回数が1とすると、第4のキーワード(1)についての共起ベクトルは、 $(0, 1, \dots, 1)$ と求まる。同様にして、他の第4のキーワード(第2のキーワード(2)、第2のキーワード(3)、・・・)についての共起ベクトルを求める。

【0309】

キーワード入力部11に入力された第2のキーワードと同一の第4のキーワードについての共起ベクトルと、対応する類似キーワードを求めたい第4のキーワードについての共起ベクトルとの類似の度合いを求める。例えば、キーワード入力部11に入力された第2のキーワードと同一の第4のキーワードについての共起ベクトルが $(a_1, a_2, a_3, \dots, a_n)$ 、対応する類似キーワードを求めたい第4のキーワードについての共起ベクトルが $(b_1, b_2, b_3, \dots, b_n)$ とすると、 $(a_1 - b_1)^2 + (a_2 - b_2)^2 + (a_3 - b_3)^2 + \dots + (a_n - b_n)^2$ の値を算出する。算出された値が両共起ベクトル間の類似の度合いを示している。算出された値が低いほど、類似の度合いが高い。

30

【0310】

算出された値が最も低いときのキーワード入力部11に入力された第2のキーワードと同一の第4のキーワードを、対応する類似キーワードを求めたい第4のキーワードに類似する類似キーワードとする。

40

【0311】

(共起ベクトルを用いる手法(2))

類似キーワード決定部41は、まず、第4のキーワードを用いて知識DB25を全文検索し、各第4のキーワードと共起して出現した語(共起語)を抽出する。そして、各第4のキーワードが、抽出された共起語と共起して知識DB25中出现した回数を要素とするベクトルを、各第4のキーワードについての共起ベクトルとして求める。

【0312】

例えば、第4のキーワード(1)が共起語 w_1 と共起して出現した回数が2、共起語 w_2 と共起して出現した回数が0、共起語 w_3 と共起して出現した回数が1・・・、パター

50

ン c_n と共起して出現した回数が 1 とすると、第 4 のキーワード (1) についての共起ベクトルは、(2 , 0 , 1 , \dots 1) と求まる。同様にして、他の第 4 のキーワード (第 2 のキーワード (2) , 第 2 のキーワード (3) , \dots) についての共起ベクトルを求める。

【 0 3 1 3 】

キーワード入力部 1 1 に入力された第 2 のキーワードと同一の第 4 のキーワードについての共起ベクトルと、対応する類似キーワードを求めたい第 4 のキーワードについての共起ベクトルとの類似の度合いを求める。例えば、キーワード入力部 1 1 に入力された第 2 のキーワードと同一の第 4 のキーワードについての共起ベクトルが (a_1 , a_2 , a_3 , \dots a_n) 、対応する類似キーワードを求めたい第 4 のキーワードについての共起ベクトルが (b_1 , b_2 , b_3 , \dots b_n) とすると、($a_1 - b_1$)² + ($a_2 - b_2$)² + ($a_3 - b_3$)² + \dots + ($a_n - b_n$)² の値を算出する。算出された値が両共起ベクトル間の類似の度合いを示している。算出された値が低いほど、類似の度合いが高い。

10

【 0 3 1 4 】

算出された値が最も低いときのキーワード入力部 1 1 に入力された第 2 のキーワードと同一の第 4 のキーワードを、対応する類似キーワードを求めたい第 4 のキーワードに類似する類似キーワードとする。

【 0 3 1 5 】

なお、本発明の実施の形態においては、類似キーワード決定部 4 1 が、知識 DB 2 5 ではなく、他の文書データを用いて上記共起ベクトルを求める構成を採ることもできる。例えば、大量の文書データが格納された大規模コーパス (図示を省略) を用いて上記共起ベクトルを求める構成を採ることもできる。

20

【 0 3 1 6 】

(シソーラスデータを用いる手法)

シソーラスデータが分類語彙表の形式で格納されているシソーラスデータベース (図 1 7 では図示を省略) を用意する。類似キーワード決定部 4 1 は、シソーラスデータベース内に格納されているシソーラスデータ中の各単語に振られた、10 桁の分類番号における各桁の数字の一致の割合を用いて、第 4 のキーワードと、キーワード入力部 1 1 に入力された第 2 のキーワード (と同一の第 4 のキーワード) との類似度を求める。

30

【 0 3 1 7 】

すなわち、例えば、分類語彙表中の、対応する類似キーワードを求めたい第 4 のキーワードと同一の単語に振られた分類番号について、キーワード入力部 1 1 に入力された第 2 のキーワード (と同一の第 4 のキーワード) と同一の単語に振られた分類番号との間での、各桁の数字の一致の割合を算出し、算出された値の大きさを類似度とする。そして、算出された値が最も大きいときの、上記第 2 のキーワード (と同一の第 4 のキーワード) を、対応する類似キーワードを求めたい第 4 のキーワードに類似する類似キーワードとして決定する。

【 0 3 1 8 】

質問作成部 4 2 は、キーワード増加部 1 2 の処理によって出力された第 3 のキーワードと第 4 のキーワードとの組み合わせによって構成される質問を作成する。

40

【 0 3 1 9 】

本発明の実施の形態においては、図 1 7 に示す構成から質問作成部 4 2 を省略し、解答候補抽出部 4 3 が、キーワード増加部 1 2 によって出力された第 3 のキーワードと第 4 のキーワードとによって構成される質問に対する解答候補を抽出し、出力する構成を採ってもよい。

【 0 3 2 0 】

解答候補抽出部 4 3 は、知識 DB 2 5 から、質問作成部 4 2 によって作成された各質問を構成する第 3 のキーワードと第 4 のキーワードを含む文書を検索し、この検索処理で抽出された文書に含まれる言語表現のうち、各質問を構成する第 4 のキーワードに類似する

50

類似キーワードに対応付けられて解答タイプ入力部 3 1 に入力された解答タイプに適合する言語表現を、解答候補として抽出する。解答候補抽出部 4 3 は、知識 DB 2 5 から、質問作成部 4 2 によって作成された各質問を構成する第 3 のキーワードと第 4 のキーワードを含む文書を検索し、この検索処理で抽出された文書に含まれる言語表現のうち、各質問を構成する第 4 のキーワードに類似する類似キーワードに予め対応付けられた解答タイプに適合する言語表現を、解答候補として抽出する構成を採ってもよい。

【0321】

また、解答候補抽出部 4 3 は、質問応答装置 4 が質問作成部 4 2 を省略する構成を採るときは、知識 DB 2 5 から、キーワード増加部 1 2 によって出力される第 3 のキーワードと第 4 のキーワードを含む文書を検索し、この検索処理で抽出された文書に含まれる言語表現のうち、第 4 のキーワードに類似する類似キーワードに対応付けられて解答タイプ入力部 3 1 に入力された（または予め類似キーワードに対応付けられた）解答タイプに適合する言語表現を、第 3 のキーワードと当該第 4 のキーワードとによって構成される各質問に対する解答候補として抽出する。

10

【0322】

すなわち、各質問に対する解答候補の抽出に、各質問を構成する第 4 のキーワードに類似する類似キーワードに対応付けられた解答タイプを用いる。

【0323】

例えば、「日本の旧首都は？」という質問に対する解答の候補を、第 4 のキーワード「旧首都」に類似する類似キーワード「首都」に対応付けられている解答タイプ「固有名詞（地名）」を用いて抽出する。

20

【0324】

解答表出力部 1 5 は、抽出された各解答候補が質問と対応付けられた表を解答表として出力する。

【0325】

図 1 8 は、本発明の第 4 の実施の形態における質問応答処理フローの一例を示す図である。キーワード入力部 1 1 に、第 1 のキーワードと第 2 のキーワードを入力キーワードとして入力する（ステップ S 3 1）。例えば、第 1 のキーワード「日本」、「アメリカ」、・・・と第 2 のキーワード「面積」、「首都」、・・・が入力される。また、解答タイプ入力部 3 1 に、第 2 のキーワードに対応付けられた解答タイプを入力する（ステップ S 3 2）。例えば、第 2 のキーワード「面積」に対応付けられた解答タイプ「固有名詞（数値）」、第 2 のキーワード「首都」に対応付けられた解答タイプ「固有名詞（地名）」が入力される。

30

【0326】

キーワード増加部 1 2 のパターン抽出部 1 2 1 で、入力キーワードをキーワード抽出用 DB 1 6 で全文検索し、複数の入力キーワードの周辺に出現したパターンを c_i として抽出する（ステップ S 3 3）。周辺に出現するパターンの定義は適宜行なう。なお、パターン c_i の抽出は、第 1 のキーワードと第 2 のキーワードそれぞれについて行う。

【0327】

キーワード増加部 1 2 のキーワード抽出部 1 2 2 で、パターン抽出部 1 2 1 で抽出したパターン c_i をキーワード抽出用 DB 1 6 で全文検索し、パターン c_i によって抽出される表現 $e x p$ を抽出すると同時に、抽出した表現 $e x p$ を $S c o r e$ の値の大きい順にソートし、キーワードとして出力する（ステップ S 3 4）。

40

【0328】

ステップ S 3 4 の処理によって、例えば、第 1 のキーワードが、多数の第 3 のキーワード（例えば、「日本」、「アメリカ」、「ドイツ」、「イタリア」、「フランス」、「イギリス」・・・）に増加する。また、第 2 のキーワードが、多数の第 4 のキーワード（例えば、「面積」、「人口」、「緯度」、「首都」、「旧首都」、「最南端都市」・・・）に増加する。

【0329】

50

類似キーワード決定部 4 1 が、第 4 のキーワードと類似する類似キーワードを決定する（ステップ S 3 5）。例えば、第 4 のキーワード「旧首都」に類似する類似キーワードとして、キーワード入力部 1 1 に入力された第 2 のキーワード（と同一の第 4 のキーワード）である「首都」が決定される。

【 0 3 3 0 】

キーワード入力部 1 1 へのキーワードの入力がある間（ステップ S 3 6）は、上述したステップ S 3 1 ~ ステップ S 3 5 の処理が繰り返される。

【 0 3 3 1 】

ステップ S 3 6 において、キーワード入力部 1 1 への入力キーワードの入力がなくなると、質問作成部 4 2 が、第 3 のキーワードと第 4 のキーワードとにより構成される質問を作成する（ステップ S 3 7）。例えば、「日本の旧首都は？」、「アメリカの面積は？」、「ドイツの緯度は？」・・・といった質問を作成する。質問応答装置 4 が質問作成部 4 2 を備えない構成を採るときは、上記ステップ S 3 7 の処理は、省略される。

【 0 3 3 2 】

解答候補抽出部 4 3 は、作成された各質問に対する解答候補を抽出する（ステップ S 3 8）。すなわち、解答候補抽出部 4 3 は、知識 DB 2 5 から、質問作成部 4 2 によって作成された各質問を構成する第 3 のキーワードと第 4 キーワードを含む文書を検索し、この検索処理で抽出された文書に含まれる言語表現のうち、各質問を構成する第 4 のキーワードが類似する類似キーワードに対応付けられた解答タイプに適合する言語表現を、解答候補として抽出する。

【 0 3 3 3 】

ステップ S 3 8 において、解答候補抽出部 4 3 は、質問応答装置 4 が質問作成部 4 2 を省略する構成を採るときは、知識 DB 2 5 から、キーワード増加部 1 2 によって出力される第 3 のキーワードと第 4 のキーワードを含む文書を検索し、この検索処理で抽出された文書に含まれる言語表現のうち、第 4 のキーワードに類似する類似キーワードに対応付けられて解答タイプ入力部 3 1 に入力された解答タイプに適合する言語表現を、第 3 のキーワードと当該第 4 のキーワードとによって構成される各質問に対する解答候補として抽出する。

【 0 3 3 4 】

例えば、第 4 のキーワード「緯度」によって構成される質問に対する解答候補の抽出には、第 4 のキーワード「緯度」が類似する類似キーワード「面積」に対応付けられた解答タイプ「固有名詞（数値）」を用いる。

【 0 3 3 5 】

また、例えば、第 4 のキーワード「旧首都」によって構成される質問に対する解答候補の抽出には、第 4 のキーワード「旧首都」が類似する類似キーワードに対応付けられた解答タイプ「固有名詞（地名）」を用いる。

【 0 3 3 6 】

そして、解答表出力部 1 5 が、解答表を出力する（ステップ S 3 9）。例えば図 1 9 に示すような解答表が出力される。

【 0 3 3 7 】

図 2 0 は、本発明の第 4 の実施の形態における質問応答装置の構成の別の例を示す図である。質問応答装置 4 0 は、入力されたキーワードを増加し、増加したキーワードにより構成される質問に対する解答を出力する装置である。図 2 0 中に示す質問応答装置 4 0 が備える構成要素のうち、図 1 に示す質問応答装置 1 が備える構成要素または図 6 に示す質問応答装置 1 0 または図 1 7 に示す質問応答装置 4 が備える構成要素と同一の符号が付けられたものは、当該質問応答装置 1 または質問応答装置 1 0 または質問応答装置 4 が備える構成要素と同様の機能を有する。本発明の実施の形態においては、図 2 0 に示す構成から質問作成部 4 2 を省略し、解答候補抽出部 4 3 が、キーワード増加部 6 0 によって出力された第 3 のキーワードと第 4 のキーワードとによって構成される質問に対する解答候補を抽出し、出力する構成を採ってもよい。

10

20

30

40

50

【 0 3 3 8 】

上記の質問応答装置 40 を用いた場合の質問応答処理フローは、図 18 に示す質問応答処理フローと、ステップ S 33、ステップ S 34 の処理が異なる以外は、同様である。質問応答装置 40 を用いた場合の質問応答処理フローの一例においては、図 18 のステップ S 33 およびステップ S 34 の代わりに、キーワード増加部 60 のキーワード抽出部 101 で、キーワード入力部 11 に入力された第 1 のキーワードと同じ分野の単語を単語データ DB 61 中から抽出し、第 3 のキーワードとして出力する。また、キーワード抽出部 101 で、キーワード入力部 11 に入力された第 2 のキーワードと同じ分野の単語を単語データ DB 61 中から抽出し、第 4 のキーワードとして出力する。

【 0 3 3 9 】

また、質問応答装置 40 を用いた場合の質問応答処理フローの別の例においては、図 18 のステップ S 33 およびステップ S 34 の代わりに、キーワード増加部 60 の類似度算出部 100 が、キーワード入力部 11 に入力された第 1 のキーワードとシソーラス DB 62 中の単語との類似度を算出し、キーワード増加部 60 のキーワード抽出部 101 が、算出された類似度が予め定めた閾値以上の単語を第 3 のキーワードとして出力する。また、類似度算出部 100 が、キーワード入力部 11 に入力された第 2 のキーワードとシソーラス DB 62 中の単語との類似度を算出し、キーワード抽出部 101 が、算出された類似度が予め定めた閾値以上の単語を第 4 のキーワードとして出力する。

【 0 3 4 0 】

また、キーワード抽出部 101 は、例えば、上記算出された類似度が大きい順に所定の個数の単語をシソーラスデータ中から取り出して、上記の第 3 のキーワード、第 4 のキーワードとして出力する構成を採ることもできる。

【 0 3 4 1 】

また、本発明の第 4 の実施の形態においては、図 17 に示す質問応答装置 4 または図 20 に示す質問応答装置 40 は、例えば、解答タイプ入力部 31 に替えて、キーワード入力部 11 に入力された第 2 のキーワードに対応付けられた疑問代名詞が入力される疑問代名詞入力部（図示を省略）と、上記疑問代名詞入力部に入力された疑問代名詞に基づいて、キーワード増加部 12（またはキーワード増加部 60）によって出力される出力キーワードによって構成される質問に対する解答の候補の言語表現の類型である解答タイプを推定する解答タイプ推定部（図示を省略）とを備える構成を採ってもよい。

【 0 3 4 2 】

上記の構成においては、解答候補抽出部 43 は、知識 DB 25 から、キーワード増加部 12（またはキーワード増加部 60）によって出力される第 3 のキーワードと第 4 のキーワードを含む文書を検索し、この検索処理で抽出された文書に含まれる言語表現のうち、第 4 のキーワードに類似する類似キーワードに対応付けられて上記疑問代名詞入力部に入力された疑問代名詞に基づいて上記解答タイプ推定部が推定した解答タイプに適合する言語表現を、第 3 のキーワードと当該第 4 のキーワードとによって構成される各質問に対する解答候補として抽出してもよい。また、解答候補抽出部 43 は、知識 DB 25 から、キーワード増加部 12（またはキーワード増加部 60）によって出力される第 3 のキーワードと第 4 のキーワードを含む文書を検索し、この検索処理で抽出された文書に含まれる言語表現のうち、第 4 のキーワードに類似する類似キーワードに対応付けられた疑問代名詞（すなわち、キーワード入力部 11 に入力された第 2 のキーワードに対応付けられるものとして予め定められた疑問代名詞）に基づいて上記解答タイプ推定部が推定した解答タイプに適合する言語表現を、第 3 のキーワードと当該第 4 のキーワードとによって構成される各質問に対する解答候補として抽出してもよい。

【 0 3 4 3 】

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

10

20

30

40

50

クを利用した送受信により提供されるものである。

【図面の簡単な説明】

【0344】

【図1】本発明の第1の実施の形態における質問応答装置の構成の一例を示す図である。

【図2】解答表の一例である。

【図3】キーワードの抽出結果に対する適合率・再現率の一例を示す図である。

【図4】正解データの一例を示す図である。

【図5】本発明の第1の実施の形態における質問応答処理フローの一例を示す図である。

【図6】本発明の第1の実施の形態における質問応答装置の構成の別の例を示す図である

。

【図7】単語データDBのデータ構成例を示す図である。

【図8】シソーラスDBのデータ構成例を示す図である。

【図9】本発明の第2の実施の形態における質問応答装置の構成の一例を示す図である。

【図10】解答表の一例である。

【図11】本発明の第2の実施の形態における質問応答処理フローの一例を示す図である

。

【図12】本発明の第2の実施の形態の変形例1の構成例を示す図である。

【図13】本発明の第3の実施の形態における質問応答装置の構成の一例を示す図である

。

【図14】本発明の第3の実施の形態における質問応答処理フローの一例を示す図である

。

【図15】解答表の一例である。

【図16】本発明の第3の実施の形態の変形例1の構成例を示す図である。

【図17】本発明の第4の実施の形態における質問応答装置の構成の一例を示す図である

。

【図18】本発明の第4の実施の形態における質問応答処理フローの一例を示す図である

。

【図19】解答表の一例である。

【図20】本発明の第4の実施の形態における質問応答装置の構成の別の例を示す図である。

【図21】解の候補と得点のリストの例である。

【図22】解の候補の得点を単純に加算する方法を用いた出力結果の例である。

【図23】質問に対する出力結果の例である。

【図24】質問に対する出力結果の例である。

【図25】質問に対する出力結果の例である。

【図26】質問に対する出力結果の例である。

【図27】サポートベクトルマシン法のマージン最大化の概念を示す図である。

【符号の説明】

【0345】

1、2、3、4、10、20、30、40 質問応答装置

11 キーワード入力部

12、18、60、63 キーワード増加部

13、23、42 質問作成部

14、24、43 解答候補抽出部

15 解答表出力部

16 キーワード抽出用DB

17 学習DB

21 疑問代名詞入力部

22 解答タイプ推定部

25 知識DB

10

20

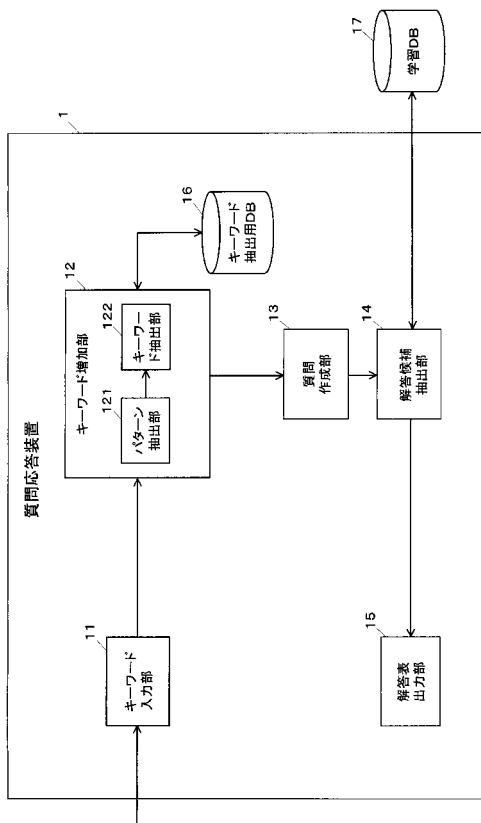
30

40

50

- 3 1 解答タイプ入力部
- 4 1 類似キーワード決定部
- 6 1 単語データDB
- 6 2 シソーラスDB
- 1 0 0 類似度算出部
- 1 0 1、1 2 2 キーワード抽出部
- 1 2 1 パターン抽出部

【図 1】



【図 2】

	面積	人口	緯度
日本	A1	A2	A3
アメリカ	B1	B2	B3
ドイツ	C1	C2	C3

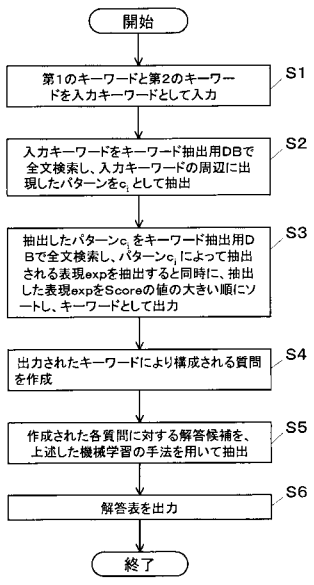
【図3】

	字種とKRを利用せず			字種の利用			KRの利用			字種とKRの利用		
	AP	RP	TP	AP	RP	TP	AP	RP	TP	AP	RP	TP
手法1	0.102	0.170	0.305	0.142	0.220	0.365	0.096	0.161	0.255	0.154	0.231	0.360
手法2	0.174	0.235	0.445	0.182	0.244	0.475	0.177	0.235	0.465	0.185	0.247	0.490
手法3	0.171	0.234	0.435	0.178	0.242	0.460	0.182	0.239	0.475	0.189	0.251	0.490
手法4	0.172	0.234	0.435	0.179	0.244	0.465	0.172	0.235	0.435	0.179	0.245	0.465
手法5	0.174	0.236	0.410	0.192	0.264	0.450	0.190	0.246	0.460	0.206	0.272	0.490

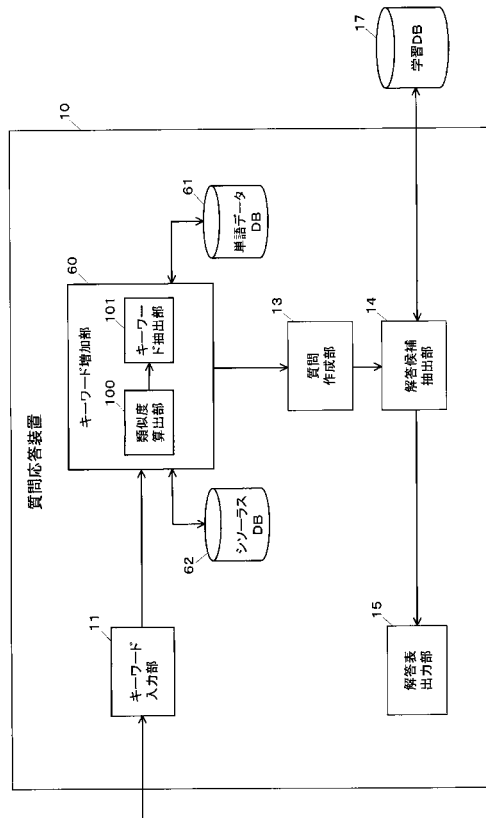
【図4】

アイスランド アイスランド共和国 ISL
 アイルランド アイルランド共和国 IRL
 アゼルバイジャン アゼルバイジャン共和国 AZE
 アゾレス諸島
 アドゥイグ アドゥイグ共和国
 アフガニスタン アフガニスタン共和国
 アメリカ アメリカ合衆国 米国 米 USA
 ...

【図5】



【図6】



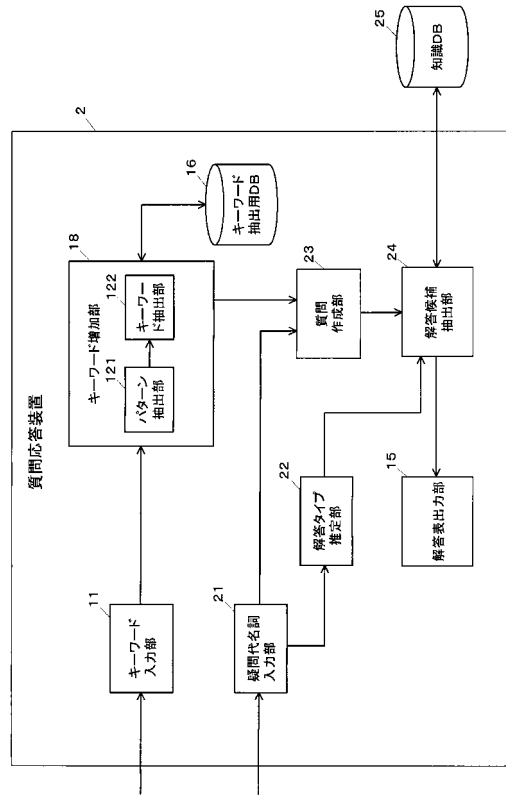
【図7】

分野	単語
国名	日本、アメリカ、ドイツ、...
数値表現	面積、人口、緯度、...
⋮	⋮

【図8】

単語	分類番号
日本	1259001012
アメリカ	1259004192
面積	1263013015
人口	1263010012
緯度	1263010015
⋮	⋮

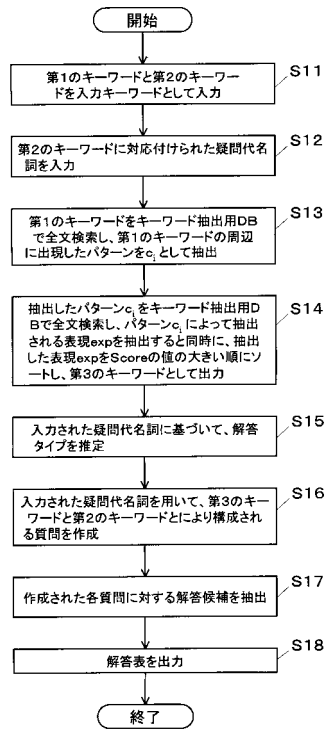
【図9】



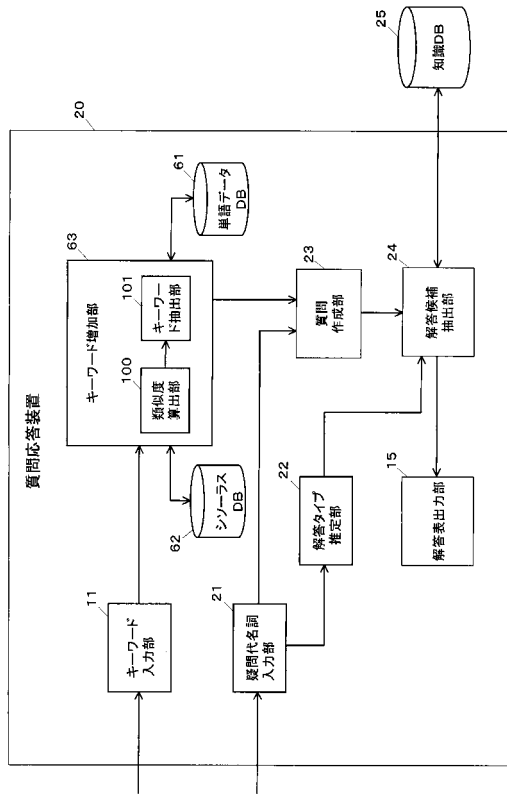
【図10】

	首都
日本	東京
アメリカ	ワシントン
ドイツ	ベルリン
⋮	⋮

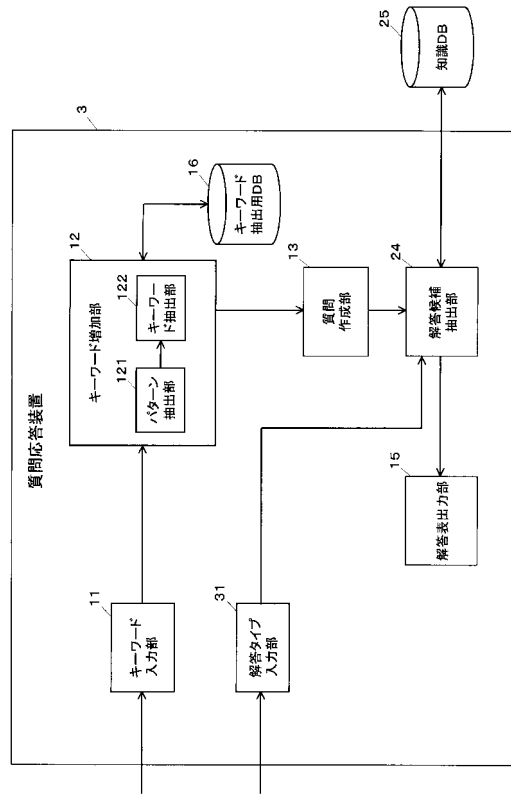
【図11】



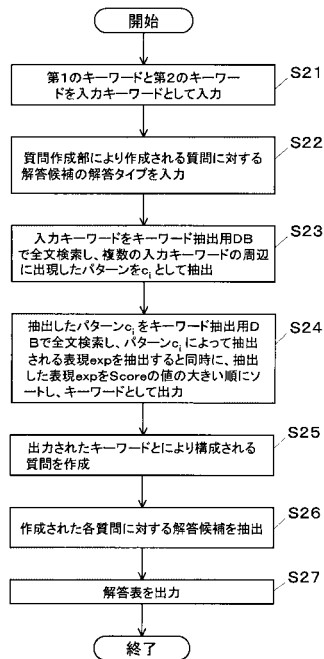
【図12】



【図13】



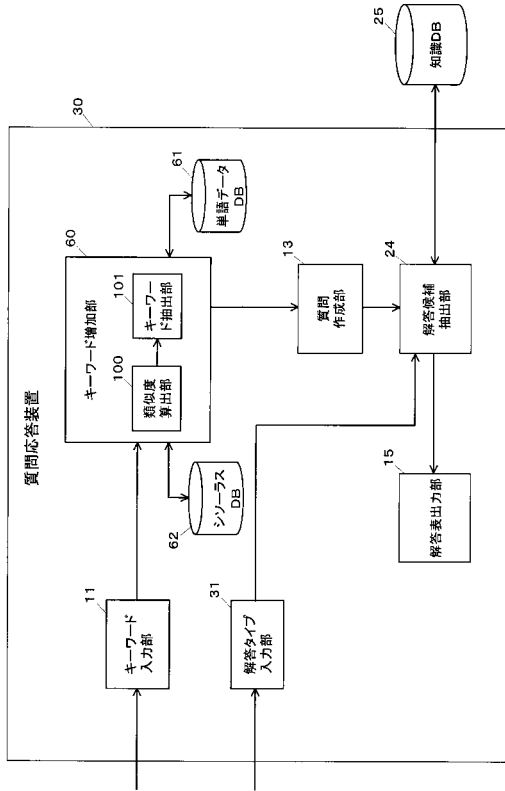
【図14】



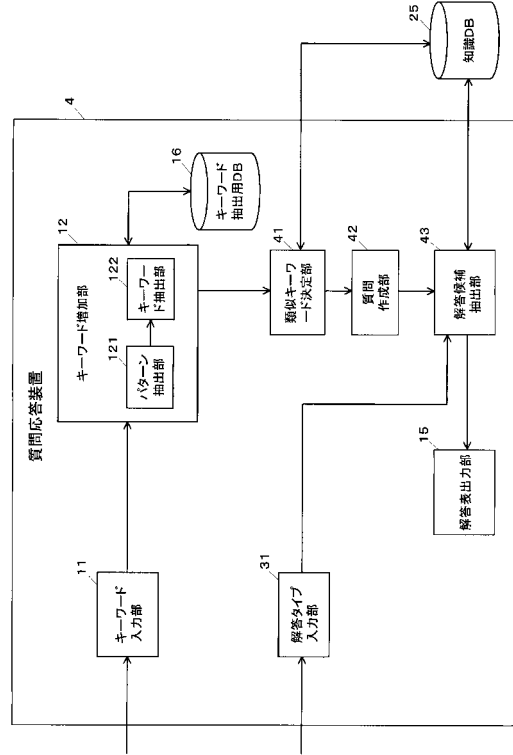
【図15】

	首都	旧首都	最南端都市
日本	東京	江戸	a
アメリカ	ワシントン	フィラデルフィア	b
ドイツ	ベルリン	ボン	c
⋮	⋮	⋮	⋮

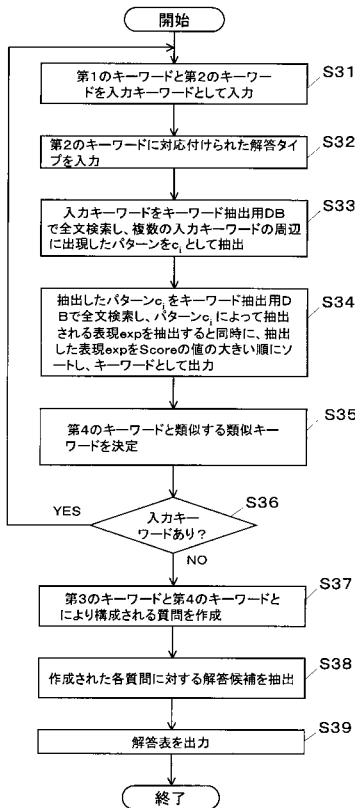
【図16】



【図17】



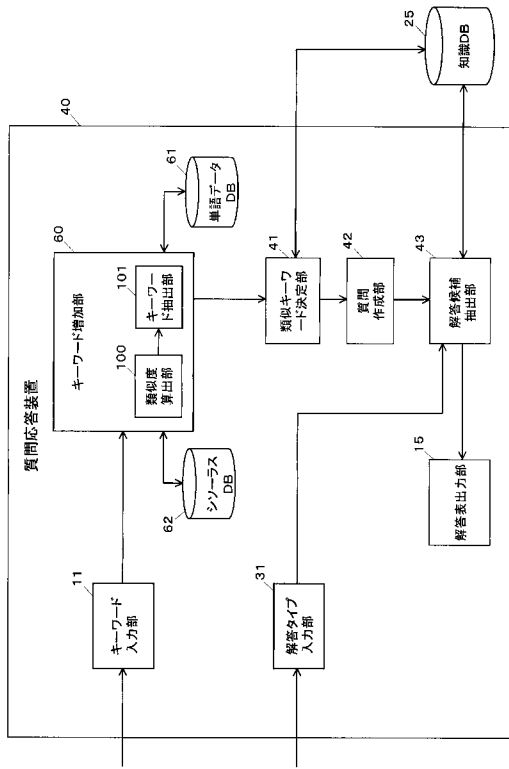
【図18】



【図19】

面積	人口	緯度	首都	旧首都	最南緯都市
日本	A1	A3	東京	江戸	a
アメリカ	B1	B3	ワシントン	フィラデルフィア	b
ドイツ	C1	C3	ベルリン	-	c
イタリア	D1	D3	ローマ	トリノ	d
フランス	E1	E3	パリ	-	e
イギリス	FC1	F3	ロンドン	-	f

【図20】



【図21】

順位	解の候補	得点	記事番号
1	京都	3.3	926324
2	東京	3.2	259312
3	東京	2.8	451245
4	東京	2.5	371922
5	東京	2.4	221328
6	北京	2.3	113127
...

【図22】

順位	解の候補	得点	記事番号
1	東京	10.9	259312, 451245, 371922, 221328
2	京都	3.3	926324
3	北京	2.3	113127
...

【図23】

順位	解の候補	得点	記事番号
1	京都	5.4	926324
2	東京	2.1	259312
3	東京	1.8	451245
4	東京	1.5	371922
5	東京	1.4	221328
6	北京	1.3	113127
...

【図25】

順位	解の候補	得点	記事番号
1	京都	5.4	926324
2	東京	2.8	259312, 451245, 371922, 221328
3	北京	1.3	113127
...

【図24】

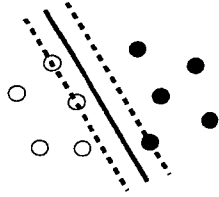
順位	解の候補	得点	記事番号
1	東京	6.8	259312, 451245, 371922, 221328
2	京都	5.4	926324
3	北京	1.3	113127
...

【図26】

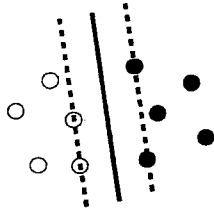
順位	解の候補	得点	記事番号
1	東京	4.3	259312, 451245, 371922, 221328
2	京都	3.3	926324
3	北京	2.3	113127
...

【 図 27 】

(A) スモールマシン



(B) ラージマシン



フロントページの続き

(72)発明者 井佐原 均

東京都小金井市貫井北町4 - 2 - 1 独立行政法人情報通信研究機構内

審査官 鈴木 和樹

(56)参考文献 特開2002 - 222147 (JP, A)

特開2005 - 157524 (JP, A)

村田真樹、外2名、質問応答システムを用いた情報抽出、言語処理学会第6回年次大会ワークショップ論文集、日本、言語処理学会、2000年 3月10日、p. 33 - 40

(58)調査した分野(Int.Cl., DB名)

G06F 17/30