

(19) 日本国特許庁(JP)

(12) 特許公報(B2)

(11) 特許番号

特許第5213098号
(P5213098)

(45) 発行日 平成25年6月19日(2013.6.19)

(24) 登録日 平成25年3月8日(2013.3.8)

(51) Int.Cl. F 1
G 0 6 F 17/30 (2006.01)
 G 0 6 F 17/30 3 3 0 C
 G 0 6 F 17/30 1 7 0 A
 G 0 6 F 17/30 2 1 0 D

請求項の数 8 外国語出願 (全 12 頁)

(21) 出願番号	特願2007-165692 (P2007-165692)	(73) 特許権者	301022471 独立行政法人情報通信研究機構 東京都小金井市貫井北町4-2-1
(22) 出願日	平成19年6月22日(2007.6.22)	(74) 代理人	100120868 弁理士 安彦 元
(65) 公開番号	特開2009-3814 (P2009-3814A)	(72) 発明者	呉 友政 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
(43) 公開日	平成21年1月8日(2009.1.8)	(72) 発明者	柏岡 秀紀 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
審査請求日	平成22年6月9日(2010.6.9)	審査官	打出 義尚

最終頁に続く

(54) 【発明の名称】 質問応答方法及びシステム

(57) 【特許請求の範囲】

【請求項1】

ユーザから入力された質問文を構成する単語から、キーワードと、質問の種別を規定するアンサータイプとを特定する質問文解析ステップと、

上記質問文解析ステップにおいて抽出したキーワードを検索クエリーとし、当該キーワードに関連するウェブページを検索する検索ステップと、

上記検索ステップにおいて検索した各ウェブページから、上記アンサータイプに基づいて解答候補を順次抽出する解答候補抽出ステップと、

上記解答候補抽出ステップにおいて抽出した各解答候補に基づいて、候補選択スニペットをクラスタに割り当て、これをトレーニングデータとする、クラスタリングステップと

10

上記トレーニングデータを解析することにより上記クラスタを分類し、更に上記トレーニングデータの解析と同一解析条件の下で上記質問文を解析し、当該質問文の解析結果と最も類似するクラスタを順に抽出し、これを応答として出力する分類ステップとを有し、

上記分類ステップは、上記トレーニングデータと上記質問文との単語重複度を示す S B F S (similarity-based feature set)、上記トレーニングデータと上記質問文とのブーリアン重複度を示す B M F S (Boolean match-based feature set)、上記トレーニングデータを構成する解答候補の前後を構成する文字を含めた文字列と上記質問文との類似度を示す W W F S (window-based word feature set) に基づいて、当該質問文の解析結果と最も類似するクラスタを順に抽出し、

20

上記 S B F S は、上記キーワードの bi-gram の一致度に基づくものであり、
 上記 B M F S は、解答候補が質問文の bi-gram と一致した bi-gram を有するか否かに基づくものであり、

上記 W W F S は、以下の I S F 値により重み付けされていること

$$I S F (w_j, C_i) = (N(w_j, C_i) + 0.5) / (N(w_j) + 0.5)$$

ここで、 $N(w_j)$ は、単語 w_j が含まれているウェブページのスニペットの総数、 $N(w_j, C_i)$ は、単語 w_j が含まれているクラスタ C_i におけるスニペットの数を特徴とする質問応答方法。

【請求項 2】

上記分類ステップは、S V M (Support Vector Machine) を利用することにより、質問文の解析結果と最も類似するクラスタを順に抽出すること
 を特徴とする請求項 1 記載の質問応答方法。

10

【請求項 3】

ユーザから入力された質問文を構成する単語から、キーワードと、質問の種別を規定するアンサータイプとを特定する質問文解析手段と、

上記質問文解析手段により抽出されたキーワードを検索クエリーとし、当該キーワードに関連するウェブページを検索する検索手段と、

上記検索ステップにより検索された各ウェブページから、上記アンサータイプに基づいて解答候補を順次抽出する解答候補抽出手段と、

上記解答候補抽出手段により抽出された各解答候補に基づいて、候補選択スニペットをクラスタに割り当て、これをトレーニングデータとする、クラスタリング手段と、
 上記トレーニングデータを解析することにより上記クラスタを分類し、更に上記トレーニングデータの解析と同一解析条件の下で上記質問文を解析し、当該質問文の解析結果と最も類似するクラスタを順に抽出し、これを応答として出力する分類手段とを備え、

20

上記分類手段は、上記トレーニングデータと上記質問文との単語重複度を示す S B F S (similarity-based feature set)、上記トレーニングデータと上記質問文とのブリアン重複度を示す B M F S (Boolean match-based feature set)、上記トレーニングデータを構成する解答候補の前後を構成する文字を含めた文字列と上記質問文との類似度を示す W W F S (window-based word feature set) に基づいて、当該質問文の解析結果と最も類似するクラスタを順に抽出し、

30

上記分類手段は、上記 S B F S を、上記キーワードの bi-gram の一致度に基づくものとし、

上記 B M F S を、解答候補が質問文の bi-gram と一致した bi-gram を有するか否かに基づくものとし、

上記 W W F S を、以下の I S F 値により重み付けすること

$$I S F (w_j, C_i) = (N(w_j, C_i) + 0.5) / (N(w_j) + 0.5)$$

ここで、 $N(w_j)$ は、単語 w_j が含まれているウェブページのスニペットの総数、 $N(w_j, C_i)$ は、単語 w_j が含まれているクラスタ C_i におけるスニペットの数を特徴とする質問応答システム。

【請求項 4】

40

上記分類手段は、S V M (Support Vector Machine) を利用することにより、質問文の解析結果と最も類似するクラスタを順に抽出すること

を特徴とする請求項 3 記載の質問応答システム。

【請求項 5】

上記質問文解析手段と、上記検索手段と、上記解答候補抽出手段と、上記クラスタリング手段と、上記分類手段とを備える制御装置と、当該制御装置に対して通信網を介して情報を送受信可能な複数のユーザ用端末装置とを備え、

上記ユーザ用端末装置は、ユーザからの上記質問文の入力を受け付け、これを通信網を介して上記制御装置における上記質問文解析手段へと送信するとともに、当該制御装置における上記分類手段から出力される上記応答を上記通信網を介して受信し、これをユーザに

50

表示すること

を特徴とする請求項 3 又は 4 記載の質問応答システム。

【請求項 6】

ユーザから入力された質問文を構成する単語から、キーワードと、質問の種別を規定するアンサータイプとを特定する質問文解析ステップと、

上記質問文解析ステップにおいて抽出したキーワードを検索クエリーとし、当該キーワードに関連するウェブページを検索する検索ステップと、

上記検索ステップにおいて検索した各ウェブページから、上記アンサータイプに基づいて解答候補を順次抽出する解答候補抽出ステップと、

上記解答候補抽出ステップにおいて抽出した各解答候補に基づいて、候補選択スニペットをクラスタに割り当て、これをトレーニングデータとする、クラスタリングステップと、

上記トレーニングデータを解析することにより上記クラスタを分類し、更に上記トレーニングデータの解析と同一解析条件の下で上記質問文を解析し、当該質問文の解析結果と最も類似するクラスタを順に抽出し、これを応答として出力する分類ステップとをコンピュータに実行させ、

10

上記分類ステップは、上記トレーニングデータと上記質問文との単語重複度を示す S B F S (similarity-based feature set)、上記トレーニングデータと上記質問文とのブリアン重複度を示す B M F S (Boolean match-based feature set)、上記トレーニングデータを構成する解答候補の前後を構成する文字を含めた文字列と上記質問文との類似度を示す W W F S (window-based word feature set) に基づいて、当該質問文の解析結果と最も類似するクラスタを順に抽出し、

20

上記 S B F S は、上記キーワードの bi-gram の一致度に基づくものであり、

上記 B M F S は、解答候補が質問文の bi-gram と一致した bi-gram を有するか否かに基づくものであり、

上記 W W F S は、以下の I S F 値により重み付けされていること

$$I S F (w_j, C_i) = (N(w_j, C_i) + 0.5) / (N(w_j) + 0.5)$$

ここで、 $N(w_j)$ は、単語 w_j が含まれているウェブページのスニペットの総数、 $N(w_j, C_i)$ は、単語 w_j が含まれているクラスタ C_i におけるスニペットの数

を特徴とするプログラム。

【請求項 7】

上記分類ステップは、S V M (Support Vector Machine) を利用することにより、質問文の解析結果と最も類似するクラスタを順に抽出すること

を特徴とする請求項 6 記載のプログラム。

30

【請求項 8】

請求項 7 項記載のプログラムが記録されていることを特徴とする記録媒体。

【発明の詳細な説明】

【技術分野】

【0001】

ユーザから入力された質問文に対して応答を出力可能な質問応答方法及びシステム、プログラム並びに記録媒体に関する。

40

【背景技術】

【0002】

近年におけるインターネットの普及に伴い、ユーザは、検索エンジンを利用して所望のウェブページをサーチし、そこから知見を得ることが可能となってきた。特にユーザが知りたい情報について検索エンジンを介してウェブページを検索する場合、検索クエリーとして、キーワードを入力することにより、当該キーワードに関連するスニペットを持つウェブページが自動抽出され、これを検索リストとして表示されることになる。ユーザは、かかる検索リストに表示されたウェブページ一覧から所望のウェブページにアクセスし、知りたい情報を取得することが可能となる。

【0003】

50

ところで、現在におけるウェブページの検索方法では、検索リストに表示されたウェブページ一覧から、所望の情報が記載されている、真のウェブページをユーザ自身が順次アクセスしながら見つけ出す必要があり、労力の負担が増大し、また検索に要する時間が長期化してしまうという問題点があった。

【0004】

このため、このような検索エンジンを介して検索リストを表示する代替として、ユーザ自身が知りたい情報を自然な文章として端末を介して入力し、かかる質問文に対する応答を直接出力する質問応答システムが従来から望まれていた。このため、かかる質問応答システムに関する研究も従来より行われていた。

【0005】

従来の質問応答システムは、4つのカテゴリーに分類することができる。

【0006】

先ず、質問文の全てのキーワードと、解答候補との間で類似性を示す距離を求め、これに基づいて解答候補から正解を選び出すモデルが提案されている。しかし、このモデルでは、質問と解答候補が依拠する文章とが表面上一致していないだけで、正解を出すことができなくなるといった問題点があった。

【0007】

また、質問文を最初に予め定義したカテゴリーに分類し、これをオフラインの下で学習したアンサーパターンを利用して正解を抽出するモデルも提案されている。しかし、このモデルは、予め定義した何種類かの質問のタイプに対しては高い正確性を出すことが可能であるが、オープンドメインな質問応答のための質問のタイプを定義するのが困難であり、あらゆるタイプの質問に対して対応することができない。

【0008】

自然言語処理(NLP: Natural Language Processing)に基づくモデルは、ユーザの質問を解析し、応答に相当する文を意味的な表現へと繋げ、そして意味的にマッチングするものを解答として見つけ出すものである。このモデルは、TREC(Text REtrieval Conference)のワークショップにおいてよく実演されるものであるが、NLPツールの高パフォーマンスに大きく依存するものである。このため、処理時間が長時間に亘るとともに、作業量の増加が無視できない。

【0009】

さらに、マシンラーニングに基づくモデルも研究されている(例えば、特許文献1参照。)が、未だその有用性は確立されていない。

【0010】

即ち、これらの開示技術は、手入力された質問と解答のペアが所定量必要になるという問題点がある。また、マシンラーニング技術に特化した上記ペアを収集するのは多大な労力が必要になり、コスト増につながる。

【非特許文献1】Jun Suzuki, Yutaka Sasaki, Eisaku Maeda. SVM Answer Selection for Open-Domain Question Answering, In Proc. of Coling-2002, pp974 ~ 980(2002).

【発明の開示】

【発明が解決しようとする課題】

【0011】

そこで、本発明は、上述した問題点に鑑みて案出されたものであり、オープンドメインな質問応答システムを実現する上で、その解答の正答率を向上させることが可能な質問応答方法及びシステム、プログラム並びに記録媒体を提供することにある。

【課題を解決するための手段】

【0012】

本発明を適用した質問応答方法は、ユーザから入力された質問文を構成する単語から、キーワードと、質問の種別を規定するアンサータイプとを特定する質問文解析ステップと、上記質問文解析ステップにおいて抽出したキーワードを検索クエリーとし、当該キーワードに関連するウェブページを検索する検索ステップと、上記検索ステップにおいて検索

10

20

30

40

50

した各ウェブページから、上記アンサータイプに基づいて解答候補を順次抽出する解答候補抽出ステップと、上記解答候補抽出ステップにおいて抽出した各解答候補に基づいて、候補選択スニペットをクラスタに割り当て、これをトレーニングデータとする、クラスタリングステップと、上記トレーニングデータを解析することにより上記クラスタを分類し、更に上記トレーニングデータの解析と同一解析条件の下で上記質問文を解析し、当該質問文の解析結果と最も類似するクラスタを順に抽出し、これを応答として出力する分類ステップとを有し、上記分類ステップは、上記トレーニングデータと上記質問文との単語重複度を示すS B F S (similarity-based feature set)、上記トレーニングデータと上記質問文とのブーリアン重複度を示すB M F S (Boolean match-based feature set)、上記トレーニングデータを構成する解答候補の前後を構成する文字を含めた文字列と上記質問文との類似度を示すW W F S (window-based word feature set)に基づいて、当該質問文の解析結果と最も類似するクラスタを順に抽出し、上記S B F Sは、上記キーワードのbi-gramの一致度に基づくものであり、上記B M F Sは、解答候補が質問文のbi-gramと一致したbi-gramを有するか否かに基づくものであり、上記W W F Sは、以下のI S F値により重み付けされていること

$$I S F (w_j, C_i) = (N(w_j, C_i) + 0.5) / (N(w_j) + 0.5)$$

ここで、 $N(w_j)$ は、単語 w_j が含まれているウェブページのスニペットの総数、 $N(w_j, C_i)$ は、単語 w_j が含まれているクラスタ C_i におけるスニペットの数、

を特徴とする。

【0013】

本発明を適用した質問応答システムは、ユーザから入力された質問文を構成する単語から、キーワードと、質問の種別を規定するアンサータイプとを特定する質問文解析手段と、上記質問文解析手段により抽出されたキーワードを検索クエリーとし、当該キーワードに関連するウェブページを検索する検索手段と、上記検索ステップにより検索された各ウェブページから、上記アンサータイプに基づいて解答候補を順次抽出する解答候補抽出手段と、上記解答候補抽出手段により抽出された各解答候補に基づいて、候補選択スニペットをクラスタに割り当て、これをトレーニングデータとする、クラスタリング手段と、上記トレーニングデータを解析することにより上記クラスタを分類し、更に上記トレーニングデータの解析と同一解析条件の下で上記質問文を解析し、当該質問文の解析結果と最も類似するクラスタを順に抽出し、これを応答として出力する分類手段とを備え、上記分類手段は、上記トレーニングデータと上記質問文との単語重複度を示すS B F S (similarity-based feature set)、上記トレーニングデータと上記質問文とのブーリアン重複度を示すB M F S (Boolean match-based feature set)、上記トレーニングデータを構成する解答候補の前後を構成する文字を含めた文字列と上記質問文との類似度を示すW W F S (window-based word feature set)に基づいて、当該質問文の解析結果と最も類似するクラスタを順に抽出し、上記分類手段は、上記S B F Sを、上記キーワードのbi-gramの一致度に基づくものとし、上記B M F Sを、解答候補が質問文のbi-gramと一致したbi-gramを有するか否かに基づくものとし、上記W W F Sを、以下のI S F値により重み付けすること

$$I S F (w_j, C_i) = (N(w_j, C_i) + 0.5) / (N(w_j) + 0.5)$$

ここで、 $N(w_j)$ は、単語 w_j が含まれているウェブページのスニペットの総数、 $N(w_j, C_i)$ は、単語 w_j が含まれているクラスタ C_i におけるスニペットの数

を特徴とする。

【0014】

本発明を適用したプログラムは、ユーザから入力された質問文を構成する単語から、キーワードと、質問の種別を規定するアンサータイプとを特定する質問文解析ステップと、上記質問文解析ステップにおいて抽出したキーワードを検索クエリーとし、当該キーワードに関連するウェブページを検索する検索ステップと、上記検索ステップにおいて検索した各ウェブページから、上記アンサータイプに基づいて解答候補を順次抽出する解答候補抽出ステップと、上記解答候補抽出ステップにおいて抽出した各解答候補に基づいて、候

10

20

30

40

50

補選択スニペットをクラスタに割り当て、これをトレーニングデータとする、クラスタリングステップと、上記トレーニングデータを解析することにより上記クラスタを分類し、更に上記トレーニングデータの解析と同一解析条件の下で上記質問文を解析し、当該質問文の解析結果と最も類似するクラスタを順に抽出し、これを応答として出力する分類ステップとをコンピュータに実行させ、上記分類ステップは、上記トレーニングデータと上記質問文との単語重複度を示すS B F S (similarity-based feature set)、上記トレーニングデータと上記質問文とのブーリアン重複度を示すB M F S (Boolean match-based feature set)、上記トレーニングデータを構成する解答候補の前後を構成する文字を含めた文字列と上記質問文との類似度を示すW W F S (window-based word feature set)に基づいて、当該質問文の解析結果と最も類似するクラスタを順に抽出し、上記S B F Sは、上記キーワードのbi-gramの一致度に基づくものであり、上記B M F Sは、解答候補が質問文のbi-gramと一致したbi-gramを有するか否かに基づくものであり、上記W W F Sは、以下のI S F値により重み付けされていること

$$I S F (w_j, C_i) = (N(w_j, C_i) + 0.5) / (N(w_j) + 0.5)$$

ここで、 $N(w_j)$ は、単語 w_j が含まれているウェブページのスニペットの総数、 $N(w_j, C_i)$ は、単語 w_j が含まれているクラスタ C_i におけるスニペットの数を特徴とする。

【発明の効果】

【0015】

上述した構成からなる本発明では、後述する表1～3の結果から示されるように、オープンドメインな質問応答システムを実現する上で、その解答の正答率を向上させることが可能となる。

【発明を実施するための最良の形態】

【0016】

以下、本発明を実施するための最良の形態として、ユーザから入力された質問文に対して応答を出力可能な質問応答システムに監視、図面を参照しながら詳細に説明をする。

【0017】

本発明を適用した質問応答システム1は、図1に示すように、質問文を入力するユーザにより操作されるユーザ端末装置11と、このユーザ端末装置11により通信網12を介してそれぞれアクセス可能なウェブサーバ13と、このウェブサーバ13を制御するための制御装置14とを備えている。

【0018】

ユーザ端末装置11は、例えばパーソナルコンピュータ(PC)等が適用され、質問文を入力するためのマウスやキーボード等からなる操作部と、情報をユーザに対して表示するための、例えば液晶ディスプレイからなる表示部を備える。このユーザ端末装置11は、ユーザからの質問文の入力を受け付けた場合に、これを通信網12を介して制御装置14へと送信する。

【0019】

通信網12は、例えばウェブサーバ13とユーザ端末装置11とを電話回線を介して接続されるインターネット網を始め、T A / モデムと接続されるI S D N (Integrated Services Digital Network) / B (broadband) - I S D N 等のように、情報の双方向送受信を可能とした公衆通信網等である。

【0020】

また制御装置14も同様にコンピュータで構成されるものであり、相互にバスで接続されたCPU (Central Processing Unit) や、メモリ、固定ディスクと、通信網12を介してユーザ端末装置11との間で情報を送受信するための通信インターフェースとを備えている。実際に、本発明に係る質問応答システム1を実行するためのプログラムは、この制御装置14における固定ディスク等にインストールされることになる。また、このプログラムは、他のCD - R O M 等に記録された記録媒体として具体化することも可能となる。

【0021】

10

20

30

40

50

制御装置 14 は、通信網 12 を介してユーザ端末装置 11 から受信した質問文を受けて、ウェブサーバ 13 へアクセスし、後述するような処理を実行することにより、上記質問文に対する応答を作り出し、通信網 12 を介してユーザ端末装置 11 へと送信する。ユーザ端末装置 11 は、送られてきた応答を液晶ディスプレイからなる表示部を介して表示する。

【0022】

次に、本発明を適用した質問応答システム 1 の動作について説明をする。

【0023】

図 2 は、質問応答システム 1 を実行する上でのフローチャートを示している。先ずステップ S1 において、ユーザからの質問文の入力を受け付ける。ちなみに、本発明は、オープンドメインの質問応答の実現を想定しているところ、ユーザは、言語や入力形式に支配されることなく、自然に質問したい内容を文章にし、これを入力していくことになる。このため、ユーザの入力すべき内容について、複雑なルールは特段存在せず、また高精度な言語解析技術も特段必要としない。

10

【0024】

この質問文は、文書検索のクエリーとなりえるキーワードの集合と、質問の種別を規程するアンサータイプから構成されることになる。例えば、「いつ潜水艦が沈んだか？」という質問文が入力された場合においてキーワードは、「潜水艦」、「沈む」に相当し、アンサータイプは、「いつ」に相当するものとなる。即ち、このキーワードは、質問文中に含まれる名詞、動詞等を規程する単語であり、アンサータイプは、時、場所、主体、数量等、実際にユーザが知りたいカテゴリを示すものである。また、ステップ S2 は、質問からキーワードを抽出し、また疑問詞に基づいて質問の解答タイプを分類する。

20

【0025】

次に、ステップ S2 へ移行し、上記ステップ S1 において抽出したキーワードを検索クエリーとし、当該キーワードに関連するウェブページを検索する。その結果、このステップ S2 においては、キーワードに関係するウェブページが順次検索されてくることになる。このステップ S2 においては、例えば検索エンジンにおけるスニペットの記載に基づいて検索を行うようにしてもよい。ちなみに、このステップ S2 において、上述した例の質問文が入力された場合には、「潜水艦」、「沈む」というキーワードに関係するウェブページが順次検索されてくることになる。

30

【0026】

次にステップ S3 へ移行し、ステップ S2 において検索した各ウェブページから、アンサータイプに基づいて解答候補を順次抽出する。即ち、アンサータイプとして、時、場所、主体、数量等の何れかがステップ S1 において抽出されているため、これに関係する解答候補を抽出してくることになる。このステップ S3 においては、例えば検索エンジンにおけるスニペットの記載から解答候補を抽出するようにしてもよい。ちなみに、このステップ S3 において、上述した例の質問文が入力された場合には、「いつ」に相当する時を表すアンサータイプに基づく解答候補を抽出してくることになる。

【0027】

次にステップ S4 へ移行し、クラスタリングを行う。このクラスタリングは、解答候補抽出ステップにおいて抽出した各解答候補に基づいて、候補選択スニペットをクラスタに割り当て、これをトレーニングデータとする。同一の解答候補を含むウェブスニペットは、同系列のスニペットとみなし、これを同一のクラスタに属させる。そして、この割り当ての結果を分類ステップのトレーニングデータとしてとして利用する。

40

【0028】

即ち、検索エンジンのトップ m 位までのスニペット $\{s_1, s_2, \dots, s_m\}$ から n 個の解答候補 $\{c_1, c_2, \dots, c_n\}$ を抽出する。これらのスニペットは、それぞれ解答候補 $\{c_i\}$ と少なくとも 1 の質問キーワード $\{q_i\}$ を保有している。そして、これらスニペット $\{s_1, s_2, \dots, s_m\}$ は、ウェブサーチ結果のクラスタリングにより、n 個のクラスタ $\{C_1, C_2, \dots, C_n\}$ へと割り当てられることになる。

50

【 0 0 2 9 】

仮にスニペットがL個の異なる解答候補を保有するものであれば、そのスニペットはL個の異なるクラスタに割り当てられることになる。また、異なるスニペットの解答候補が互いに同一であれば、これらのスニペットは同一のクラスタに割り当てられることになる。

【 0 0 3 0 】

最終的に、クラスタ $\{C_i\}$ は、解答候補 $\{c_i\}$ の数によって決定される。そしてクラスタ C_i のクラスタ名は、解答候補 c_i に基づくものとなる。これらクラスタ化された解答候補 c_i が上述したトレーニングデータとなる。

【 0 0 3 1 】

次に、ステップS5へ移行し、トレーニングデータを解析することにより上記クラスタを分類する。

【 0 0 3 2 】

さらに次にステップS6へ移行し、トレーニングデータの解析と同一解析条件の下で上記質問文を解析し、当該質問文の解析結果と最も類似するクラスタを順に抽出し、これを応答として出力する。この分類ステップS6は、ユーザの質問から分析したテストベクトルを利用することにより、クラスタの一つに割り当て、質問に対する解答を、質問のクラスタの名称と仮定する。

【 0 0 3 3 】

以下、このステップS5とステップS6を一つの分類ステップとして説明をしていく。この分類ステップでは、SVM(Support Vector Machine)を利用することにより、質問文の解析結果と最も類似するクラスタを順に抽出するようにしてもよい。

【 0 0 3 4 】

また、この分類ステップでは、トレーニングデータと質問文との単語重複度を示すSBFS(similarity-based feature set)、トレーニングデータと質問文とのブーリアン重複度を示すBMFS(Boolean match-based feature set)、トレーニングデータを構成する解答候補の前後を構成する文字を含めた文字列と上記質問文との類似度を示すWWFS(window-based word feature set)に基づいて、当該質問文の解析結果と最も類似するクラスタを順に抽出するようにしてもよい。

【 0 0 3 5 】

SBFSとしては、キーワードの重複度、キーワードの非重複度、キーワードのbi-gramの一致度、シーソラスの重複度、キーワードと解答候補との間の規格化距離の何れか1以上に基づくものであってもよい。

【 0 0 3 6 】

BMFSは、人名が一致しているか否か、地域名が一致しているか否か、組織名が一致しているか否か、時を示す単語が一致しているか否か、数量を示す単語が一致しているか否か、語源が一致しているか否か、解答候補が質問文のbi-gramと一致したbi-gramを有するか否か、解答候補がネームディエンティタイプを要求されているか否かの何れか1以上に基づくものであってもよい。

【 0 0 3 7 】

WWFSは、以下のISF値により重み付けされていてもよい。

$$ISF(w_j, C_i) = (N(w_j, C_i) + 0.5) / (N(w_j) + 0.5)$$

ここで、 $N(w_j)$ は、単語 w_j が含まれているウェブページのスニペットの総数、 $N(w_j, C_i)$ は、単語 w_j が含まれているクラスタ C_i におけるスニペットの数である。

【 0 0 3 8 】

最後にステップS7に移行し、実際に上記プロセスの下で得た回答をユーザ端末装置11の表示部を介して表示する。

【 0 0 3 9 】

次に、本発明を用いた質問応答システム1による効果について説明をする。

【 0 0 4 0 】

10

20

30

40

50

中国語によるウェブの質問応答 (QA) における 3 種類のデータセットとしての CTREC04、CTREC05、CTEST05 を利用し、本発明を適用した U (unsupervised) -SVM を検証することとした。CTREC04 は、TREC2004FACTOID のテスト問題から翻訳された 178 個の中国語による質問のセットからなる。CTREC05 は、TREC2005FACTOID のテスト問題から翻訳された 279 個の中国語による質問のセットからなる。CTEST05 は、中国語で記載されたものを除く、TREC のテスト問題に類似する 178 個の中国語の質問のセットからなる。

【0041】

実験は、3 つの評価項目、即ち、top_1, top_5, mrr_5 に基づいて評価を行った。評価結果を表 1、2、3 に示す。ここで、top_1 は、解答の正確性がトップ 1 位である解答が含まれている割合を示している。top_5 は、解答の正確性がトップ 5 位以内である解答の一つが含まれている割合を示している。mrr_5 は、各質問に対する正解の平均相対ランク ($1/n$) を示しており、ここで最高ランク n ($n \leq 5$) としている。

10

【0042】

各データセット (CTREC04、CTREC05、CTEST05) を U-SVM を用いて解析することにより得られた各評価項目 (top_1, top_5, mrr_5) を表 1 に示す。

【0043】

【表 1】

	TREC04	TREC05	Test05
top_1	60.82%	57.33%	59.09%
mrr_5	71.31%	65.61%	67.34%
top_5	88.66%	80.00%	81.82%

20

【0044】

また、表 2 において、CTREC04 と、CTREC05 のテストデータを U-SVM と the Retrieval-M (従来の検索手法) それぞれを用いて解析することにより得られた各評価項目 (top_1, top_5, mrr_5) の相対比較を表 2 に示す。

【0045】

【表 2】

		Retrieval-M	U-SVM
Ctrec04	top_1	27.84%	53.61%
	mrr_5	43.67%	66.25%
	top_5	71.13%	88.66%
Ctrec05	top_1	34.00%	50.00%
	mrr_5	48.20%	62.38%
	top_5	71.33%	82.67%

40

【0046】

さらに、Pattern-M (パターン重視による手法) と S-SVM (SVM を利用した教化学習法によるもの) に対する U-SVM のパフォーマンス性を比較するために、CTEST05 のデータセットを用いて検証を行った。表 3 は、U-SVM、Pattern-M、S-SVM の各モデルを用いて CTEST05 を解析することにより得られた各評価項目 (top_1, mrr_5) の相対比較を示している。

【0047】

【表 3】

	S-SVM	Pattern-M	U-SVM
top_1	44.89%	53.14%	59.09%
mrr_5	56.49%	61.28%	67.34%
top_5	74.43%	73.14%	81.82%

10

【0048】

上述した表 1 ~ 3 の結果から、各モデルによる正答率（パフォーマンスランキング）は、U-SVM > Pattern-M > S-SVM > Retrieval-M の順となった。

【0049】

即ち、本発明では、上述した図 2 に示すフローに基づいて、解答を抽出していくため、正答率を向上させることが可能となる。

【図面の簡単な説明】

【0050】

【図 1】本発明を適用した質問応答システムの構成例を示す図である。

【図 2】本発明を適用した質問応答システムの処理手順を示すフローチャートである。

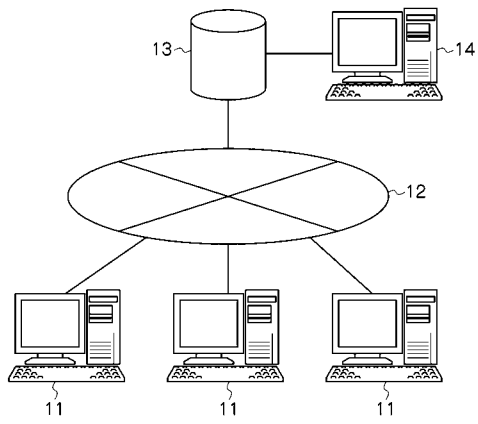
20

【符号の説明】

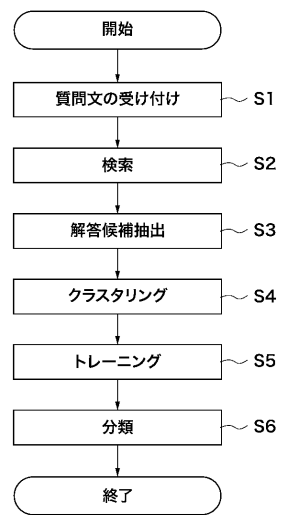
【0051】

- 1 質問応答システム
- 1 1 ユーザ端末装置
- 1 2 通信網
- 1 3 ウェブサーバ
- 1 4 制御装置

【図1】



【図2】



フロントページの続き

- (56)参考文献 特開2006-244102(JP,A)
特開平09-231238(JP,A)
特開2003-150624(JP,A)
永田昌明、外2名、日本語自然文検索システム Web Answers, 言語処理学会第12
回年次大会発表論文集, 言語処理学会, 2006年 3月13日, p.320-323
佐々木裕, SVMを用いた学習型質問応答システムSAIQA-II, 情報処理学会論文誌, 日
本, 社団法人情報処理学会, 2004年 2月15日, 第45巻, 第2号, pp.635-64
6

- (58)調査した分野(Int.Cl., DB名)
G06F 17/30