

(19) 日本国特許庁(JP)

(12) 公開特許公報(A)

(11) 特許出願公開番号

特開2014-112316

(P2014-112316A)

(43) 公開日 平成26年6月19日(2014.6.19)

(51) Int.Cl.

G06F 17/30 (2006.01)

F I

G06F 17/30 180A

G06F 17/30 210D

テーマコード (参考)

審査請求 未請求 請求項の数 10 O L (全 17 頁)

(21) 出願番号 特願2012-266589 (P2012-266589)
 (22) 出願日 平成24年12月5日 (2012.12.5)

(71) 出願人 000208891
 K D D I 株式会社
 東京都新宿区西新宿二丁目3番2号
 (74) 代理人 100135068
 弁理士 早原 茂樹
 (72) 発明者 松本 一則
 埼玉県ふじみ野市大原二丁目1番15号
 株式会社K D D I 研究所内
 (72) 発明者 服部 元
 埼玉県ふじみ野市大原二丁目1番15号
 株式会社K D D I 研究所内
 (72) 発明者 小野 智弘
 埼玉県ふじみ野市大原二丁目1番15号
 株式会社K D D I 研究所内

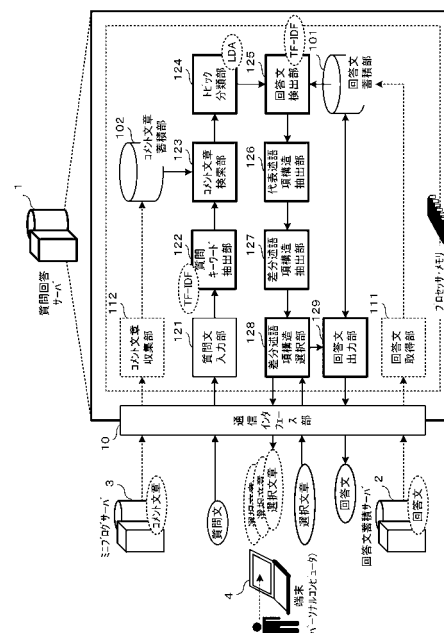
(54) 【発明の名称】 大量のコメント文章を用いた質問回答プログラム、サーバ及び方法

(57) 【要約】

【課題】ユーザの質問文に対して複数の回答文の候補が存在する場合、ユーザの意図を反映した回答文を明示する（に絞り込む）ことができる質問回答サーバ等を提供する。

【解決手段】質問文から複数の質問キーワードを抽出する質問キーワード抽出手段と、質問キーワードを含むコメント文章を検索するコメント文章検索手段と、検索された複数のコメント文章を、複数のトピックグループに分類するトピック分類手段と、各トピックと各回答文との類似度が所定閾値以上となる回答文を対応付ける回答文検出手段と、各トピックについて代表述語項構造を抽出する代表述語項構造抽出手段と、当該トピックのみに出現する代表述語項構造を、差分述語項構造として抽出する差分述語項構造抽出手段と、複数の差分述語項構造をユーザに明示し且つ選択させる差分述語項構造選択手段と、その差分述語項構造に対応する回答文をユーザに明示する回答文出力手段とを有する。

【選択図】図2



【特許請求の範囲】**【請求項 1】**

多数のコメント文章を蓄積したコメント文章蓄積部と、多数の回答文を蓄積した回答文蓄積部とを有し、ユーザからの質問文に対する回答文を抽出するようにコンピュータを機能させる質問回答プログラムであって、

質問文を入力する質問文入力手段と、

前記質問文に含まれる複数の質問キーワードを抽出する質問キーワード抽出手段と、

前記コメント文章蓄積部を用いて、前記質問キーワードを含むコメント文章を検索するコメント文章検索手段と、

検索された複数のコメント文章を、述語項構造解析によって、述語項構造の分布から複数のトピックグループに分類するトピック分類手段と、

各トピックグループに含まれるコメント文章群と、各回答文に含まれる文章との間の類似度を算出し、各トピックグループに前記類似度が所定閾値以上となる回答文を対応付ける回答文検出手段と、

各トピックグループについて、対応付けられた回答文に含まれる述語項構造の中で、当該トピックグループを特徴付ける代表述語項構造を抽出する代表述語項構造抽出手段と、

各トピックグループについて、当該トピックグループのみに出現する代表述語項構造を、差分述語項構造として抽出する差分述語項構造抽出手段と、

前記回答文検出手段によって検出された前記回答文を、対応する 1 つ以上の前記差分述語項構造に基づく文章と共に明示する回答文出力手段と

してコンピュータを機能させることを特徴とする質問回答プログラム。

【請求項 2】

複数の前記差分述語項構造に基づく文章を、ユーザインタフェースを介してユーザに明示すると共に、ユーザ操作に応じていずれか 1 つの差分述語項構造を選択させる差分述語項構造選択手段を更に有し、

前記回答文出力手段は、選択された文章の差分述語項構造に対応する回答文を、ユーザインタフェースを介して明示する

ようにコンピュータを機能させることを特徴とする請求項 1 に記載の質問回答プログラム。

【請求項 3】

前記トピック分類手段は、当該コメント文章を、分類された各トピックグループに属する確からしさ(トピック比率)を算出する LDA (Latent Dirichlet Allocation) アルゴリズムを用いて、いずれか 1 つのトピックグループに分類するようにコンピュータを機能させることを特徴とする請求項 1 又は 2 に記載の質問回答プログラム。

【請求項 4】

前記回答文検出手段は、

各トピックグループに含まれるコメント文章群から、述語項構造解析によって述語項構造を抽出すると共に、前記トピックグループにおける第 1 の特徴ベクトルを算出し、

前記回答文蓄積部に蓄積された各回答文から、述語項構造解析によって述語項構造を抽出すると共に当該回答文における第 2 の特徴ベクトルとを算出し、

前記トピックグループの第 1 のベクトルと、前記回答文の第 2 のベクトルとの間のコサイン距離に基づいて類似度を算出する

ようにコンピュータを機能させることを特徴とする請求項 1 から 3 のいずれか 1 項に記載の質問回答プログラム。

【請求項 5】

前記代表述語項構造抽出手段は、各トピックグループの代表述語項構造を、赤池情報量基準に応じて優先順に並べるようにコンピュータを機能させることを特徴とする請求項 1 から 4 のいずれか 1 項に記載の質問回答プログラム。

【請求項 6】

前記コメント文章は、不特定多数の第三者によって投稿されたものであって、

前記コメント文章蓄積部は、ミニブログ(mini Web log)サーバに投稿されたコメント文章を収集し蓄積したものであるようにコンピュータを機能させることを特徴とする請求項 1 から 5 のいずれか 1 項に記載の質問回答プログラム。

【請求項 7】

多数のコメント文章を蓄積したコメント文章蓄積部と、多数の回答文を蓄積した回答文蓄積部とを有し、ユーザからの質問文に対する回答文を抽出する質問回答サーバであって、

端末から、質問文を入力する質問文入力手段と、

前記質問文に含まれる複数の質問キーワードを抽出する質問キーワード抽出手段と、

前記コメント文章蓄積部を用いて、前記質問キーワードを含むコメント文章を検索するコメント文章検索手段と、

検索された複数のコメント文章を、述語項構造解析によって、述語項構造の分布から複数のトピックグループに分類するトピック分類手段と、

各トピックグループに含まれるコメント文章群と、各回答文との間の類似度を算出し、各トピックグループに前記類似度が所定閾値以上となる回答文を対応付ける回答文検出手段と、

各トピックグループについて、対応付けられた回答文に含まれる述語項構造の中で、当該トピックグループを特徴付ける代表述語項構造を抽出する代表述語項構造抽出手段と、

各トピックグループについて、当該トピックグループのみに出現する代表述語項構造を、差分述語項構造として抽出する差分述語項構造抽出手段と、

前記回答文検出手段によって検出された前記回答文を、対応する 1 つ以上の前記差分述語項構造に基づく文章と共に明示する回答文出力手段と

を有することを特徴とする質問回答サーバ。

【請求項 8】

複数の前記差分述語項構造に基づく文章を、ユーザインタフェースを介してユーザに明示すると共に、ユーザ操作に応じていずれか 1 つの差分述語項構造を選択させる差分述語項構造選択手段を更に有し、

前記回答文出力手段は、選択された文章の差分述語項構造に対応する回答文を、ユーザインタフェースを介して明示する

ことを特徴とする請求項 7 に記載の質問回答サーバ。

【請求項 9】

多数のコメント文章を蓄積したコメント文章蓄積部と、多数の回答文を蓄積した回答文蓄積部とを有し、ユーザからの質問文に対する回答文を抽出する装置における質問回答方法であって、

質問文を入力する第 1 のステップと、

前記質問文に含まれる複数の質問キーワードを抽出する第 2 のステップと、

前記コメント文章蓄積部を用いて、前記質問キーワードを含むコメント文章を検索する第 3 のステップと、

検索された複数のコメント文章を、述語項構造解析によって、述語項構造の分布から複数のトピックグループに分類する第 4 のステップと、

各トピックグループに含まれるコメント文章群と、各回答文との間の類似度を算出し、各トピックグループに前記類似度が所定閾値以上となる回答文を対応付ける第 5 のステップと、

各トピックグループについて、対応付けられた回答文に含まれる述語項構造の中で、当該トピックグループを特徴付ける代表述語項構造を抽出する第 6 のステップと、

各トピックグループについて、当該トピックグループのみに出現する代表述語項構造を、差分述語項構造として抽出する第 7 のステップと、

第 5 のステップによって検出された前記回答文を、対応する 1 つ以上の前記差分述語項構造に基づく文章と共に明示する第 8 のステップと

を有することを特徴とする質問回答方法。

【請求項 10】

第 8 のステップについて、

複数の前記差分述語項構造に基づく文章を、ユーザインタフェースを介してユーザに明示すると共に、ユーザ操作に応じていずれか 1 つの差分述語項構造を選択させ、

選択された文章の差分述語項構造に対応する回答文を、ユーザインタフェースを介して明示することを特徴とする請求項 9 に記載の質問回答方法。

【発明の詳細な説明】

【技術分野】

【0001】

本発明は、質問文の入力に対して最適な回答文を出力する質問回答プログラムの技術に関する。

【背景技術】

【0002】

近年、FAQ (Frequently Asked Questions) に基づく質問回答システムが構築されている。「FAQ」とは、多数の人が共通して頻繁に尋ねる質問に対する回答をまとめた問答集をいう。質問回答システムは、特定種類の情報に関する質問文をユーザから自然言語で入力し、その回答文を出力するソフトウェアをいう。一般に、質問回答システムは、仮想質問文とそれに紐づけられた回答候補文とを予めデータベースに記憶する。その上で、質問回答システムは、以下のようなステップで処理を実行する。

(1) ユーザから入力された質問文から、特徴的な単語をクエリとして抽出する。

(2) 検索エンジンを用いて、複数のクエリの出現頻度が高い仮想質問文を選択する。

(3) 選択された仮想選択文に対する回答文を選択する。

(4) 選択された回答文をユーザに提示する。

【0003】

このような質問回答システムは、ユーザに対して単体装置として存在するものもあれば、インターネット上に質問回答サーバとして接続されたものもある。この質問回答サーバは、ユーザ操作の端末からネットワークを介して質問文を受信し、回答文をその端末へ送信する(例えば非特許文献 1 参照)。

【0004】

また、他の技術として、インターネット上に、ブログ(Web log)サーバやミニブログ(mini Web log)(例えばtwitter(登録商標))サーバが接続されている。このようなブログサーバは、不特定多数の第三者からのコメント文章を受信し、他の第三者へ公開する。このようなコメント文章は、様々な話題について公開されており、勿論、前述した質問回答システムに入出力される質問文及び回答文に関連するコメント文章も多く議論されている。

【先行技術文献】

【特許文献】

【0005】

【特許文献 1】特開 2011-81626 号公報

【特許文献 2】特開 2005-141428 号公報

【特許文献 3】特開 2005-284209 号公報

【非特許文献】

【0006】

【非特許文献 1】KDDI、「au one NET コンシェルジュ」、[online]、[平成 24 年 10 月 7 日検索]、インターネット<URL: http://conciierge.auone-net.jp/inagoNetPeople/BrowserClient/GUI/kddi_missConcie3/help/help.html>

【非特許文献 2】坪坂正志、「Latent Dirichlet Allocation 入門」、[online]、[平成 24 年 10 月 7 日検索]、インターネット<URL: <http://www.slideshare.net/tsubosaka/tokyotextmining>>

【非特許文献 3】榎博史、松本一則、黒岩真吾、橋本和夫、「再起演算を用いた自然言語

10

20

30

40

50

変換方式」、電子情報通信学会論文誌(D-II), Vol.J72-D-II, No.12, pp.2080-2093, Dec. 1989、[online]、[平成24年10月7日検索]、インターネット<URL:http://jglobal.jst.go.jp/public/20090422/200902065122383276>

【発明の概要】

【発明が解決しようとする課題】

【0007】

しかしながら、同じ質問文であっても、そのユーザの質問の意図が複数あり得る場合がある。このような場合、ユーザに対して、適切な回答文が返答されない場合が多い。

【0008】

ユーザの質問文の例

Q「携帯電話機の紛失」

この質問文に対して、質問回答システムは、以下の2つキーワードを抽出する。

「携帯電話機」「紛失」

これらキーワードをクエリとして回答文を検索すると、複数の回答の選択肢がある。

A「携帯探せて安心サービスの申込方法」に関する回答文

A「携帯探せて安心サービスの利用方法」に関する回答文

この場合、ユーザとしては、紛失した携帯電話機を遠隔からロックする「利用方法」を問い合わせたつもりであるにも拘わらず、質問回答システムは、「申込方法」について回答してしまう場合もある。

【0009】

そこで、本発明は、ユーザの質問文に対して複数の回答文の候補が存在する場合、ユーザの意図を反映した回答文を明示する（に絞り込む）ことができる質問回答プログラム、サーバ及び方法を提供することを目的とする。

【課題を解決するための手段】

【0010】

本発明によれば、多数のコメント文章を蓄積したコメント文章蓄積部と、多数の回答文を蓄積した回答文蓄積部とを有し、ユーザからの質問文に対する回答文を抽出するようにコンピュータを機能させる質問回答プログラムであって、

質問文を入力する質問文入力手段と、

質問文に含まれる複数の質問キーワードを抽出する質問キーワード抽出手段と、

コメント文章蓄積部を用いて、質問キーワードを含むコメント文章を検索するコメント文章検索手段と、

検索された複数のコメント文章を、述語項構造解析によって、述語項構造の分布から複数のトピックグループに分類するトピック分類手段と、

各トピックグループに含まれるコメント文章群と、各回答文に含まれる文章との間の類似度を算出し、各トピックグループに類似度が所定閾値以上となる回答文を対応付ける回答文検出手段と、

各トピックグループについて、対応付けられた回答文に含まれる述語項構造の中で、当該トピックグループを特徴付ける代表述語項構造を抽出する代表述語項構造抽出手段と、

各トピックグループについて、当該トピックグループのみに出現する代表述語項構造を、差分述語項構造として抽出する差分述語項構造抽出手段と、

回答文検出手段によって検出された回答文を、対応する1つ以上の差分述語項構造に基づく文章と共に明示する回答文出力手段と

してコンピュータを機能させることを特徴とする。

【0011】

本発明の質問回答プログラムにおける他の実施形態によれば、

複数の差分述語項構造に基づく文章を、ユーザインタフェースを介してユーザに明示すると共に、ユーザ操作に応じていずれか1つの差分述語項構造を選択させる差分述語項構造選択手段を更に有し、

回答文出力手段は、選択された文章の差分述語項構造に対応する回答文を、ユーザイン

10

20

30

40

50

タフェースを介して明示する
ようにコンピュータを機能させることも好ましい。

【0012】

本発明の質問回答プログラムにおける他の実施形態によれば、トピック分類手段は、当該コメント文章を、分類された各トピックグループに属する確からしさ（トピック比率）を算出するLDA (Latent Dirichlet Allocation) アルゴリズムを用いて、いずれか1つのトピックグループに分類するようにコンピュータを機能させることも好ましい。

【0013】

本発明の質問回答プログラムにおける他の実施形態によれば、
回答文検出手段は、

10

各トピックグループに含まれるコメント文章群から、述語項構造解析によって述語項構造を抽出すると共に、トピックグループにおける第1の特徴ベクトルを算出し、

回答文蓄積部に蓄積された各回答文から、述語項構造解析によって述語項構造を抽出すると共に当該回答文における第2の特徴ベクトルとを算出し、

トピックグループの第1のベクトルと、回答文の第2のベクトルとの間のコサイン距離に基づいて類似度を算出する

ようにコンピュータを機能させることも好ましい。

【0014】

本発明の質問回答プログラムにおける他の実施形態によれば、代表述語項構造抽出手段は、各トピックグループの代表述語項構造を、赤池情報量基準に応じて優先順に並べるようにコンピュータを機能させることも好ましい。

20

【0015】

本発明の質問回答プログラムにおける他の実施形態によれば、

コメント文章は、不特定多数の第三者によって投稿されたものであって、

コメント文章蓄積部は、ミニブログ(mini Web log)サーバに投稿されたコメント文章を収集し蓄積したものであるようにコンピュータを機能させることも好ましい。

【0016】

本発明によれば、多数のコメント文章を蓄積したコメント文章蓄積部と、多数の回答文を蓄積した回答文蓄積部とを有し、ユーザからの質問文に対する回答文を抽出する質問回答サーバであって、

30

端末から、質問文を入力する質問文入力手段と、

質問文に含まれる複数の質問キーワードを抽出する質問キーワード抽出手段と、

コメント文章蓄積部を用いて、質問キーワードを含むコメント文章を検索するコメント文章検索手段と、

検索された複数のコメント文章を、述語項構造解析によって、述語項構造の分布から複数個のトピックグループに分類するトピック分類手段と、

各トピックグループに含まれるコメント文章群と、各回答文との間の類似度を算出し、各トピックグループに類似度が所定閾値以上となる回答文を対応付ける回答文検出手段と

、

各トピックグループについて、対応付けられた回答文に含まれる述語項構造の中で、当該トピックグループを特徴付ける代表述語項構造を抽出する代表述語項構造抽出手段と、

40

各トピックグループについて、当該トピックグループのみに出現する代表述語項構造を、差分述語項構造として抽出する差分述語項構造抽出手段と、

回答文検出手段によって検出された回答文を、対応する1つ以上の差分述語項構造に基づく文章と共に明示する回答文出力手段と

を有することを特徴とする。

【0017】

本発明の質問回答サーバにおける他の実施形態によれば、

複数の差分述語項構造に基づく文章を、ユーザインタフェースを介してユーザに明示すると共に、ユーザ操作に応じていずれか1つの差分述語項構造を選択させる差分述語項構

50

造選択手段を更に有し、

回答文出力手段は、選択された文章の差分述語項構造に対応する回答文を、ユーザインタフェースを介して明示することも好ましい。

【0018】

本発明によれば、多数のコメント文章を蓄積したコメント文章蓄積部と、多数の回答文を蓄積した回答文蓄積部とを有し、ユーザからの質問文に対する回答文を抽出する装置における質問回答方法であって、

質問文を入力する第1のステップと、

質問文に含まれる複数の質問キーワードを抽出する第2のステップと、

コメント文章蓄積部を用いて、質問キーワードを含むコメント文章を検索する第3のステップと、

検索された複数のコメント文章を、述語項構造解析によって、述語項構造の分布から複数のトピックグループに分類する第4のステップと、

各トピックグループに含まれるコメント文章群と、各回答文との間の類似度を算出し、各トピックグループに類似度が所定閾値以上となる回答文を対応付ける第5のステップと、

各トピックグループについて、対応付けられた回答文に含まれる述語項構造の中で、当該トピックグループを特徴付ける代表述語項構造を抽出する第6のステップと、

各トピックグループについて、当該トピックグループのみに出現する代表述語項構造を、差分述語項構造として抽出する第7のステップと、

第5のステップによって検出された回答文を、対応する1つ以上の差分述語項構造に基づく文章と共に明示する第8のステップとを有することを特徴とする。

【0019】

本発明の質問回答方法における他の実施形態によれば、

第8のステップについて、

複数の差分述語項構造に基づく文章を、ユーザインタフェースを介してユーザに明示すると共に、ユーザ操作に応じていずれか1つの差分述語項構造を選択させ、

選択された文章の差分述語項構造に対応する回答文を、ユーザインタフェースを介して明示することも好ましい。

【発明の効果】

【0020】

本発明の質問回答プログラム、サーバ及び方法によれば、ユーザの質問文に対して複数の回答文の候補が存在する場合、ユーザの意図を反映した回答文を明示する（に絞り込む）ことができる。

【図面の簡単な説明】

【0021】

【図1】本発明におけるシステム構成図である。

【図2】本発明における質問回答サーバの機能構成図である。

【図3】質問キーワード抽出部及びコメント文章検索部の処理を表す説明図である。

【図4】トピック分類部の処理を表す説明図である。

【図5】回答文検出部の処理を表す説明図である。

【図6】代表述語項構造抽出部、差分述語項構造抽出部、差分述語項構造選択部及び回答文出力部の処理を表す説明図である。

【図7】本発明におけるシーケンス図である。

【発明を実施するための形態】

【0022】

以下、本発明の実施の形態について、図面を用いて詳細に説明する。

【0023】

図1は、本発明におけるシステム構成図である。

【 0 0 2 4 】

図 1 によれば、インターネット上に、本発明における質問回答サーバ 1 が接続されている。質問回答サーバ 1 は、回答文を予め蓄積しているものであってもよいし、他の回答文蓄積サーバ 2 から回答文を受信するものであってもよい。尚、本発明によれば、FAQ のような質問文候補と回答文候補とを予め紐付けて記憶しておく必要はない。あくまで、回答文候補のみを予め蓄積している。

【 0 0 2 5 】

質問者が操作する端末 4 は、アクセスネットワーク及びインターネットを介して、質問回答サーバ 1 へアクセスする。そして、端末 4 は、質問文を質問回答サーバ 1 へ送信し、これに対し、質問回答サーバ 1 から回答文を受信する。以下の実施形態の中では、質問者が自然言語のテキストで端末 4 へ入力することを想定しているが、質問者が音声で入力しテキストに変換されたものであってもよい。

【 0 0 2 6 】

また、図 1 によれば、不特定多数の第三者から投稿されたコメント文章を公開するブログサーバ 3 が、インターネットに更に接続されている。ブログサーバ 3 は、例えば twitter (登録商標) サーバのようなミニブログサーバである。不特定多数の第三者は、自ら所持する端末 5 を用いて、ミニブログサーバ 3 へコメント文章を自由に投稿することができる。

【 0 0 2 7 】

本発明における質問回答サーバ 1 は、ミニブログサーバ 3 から大量のコメント文章を収集する。そして、質問回答サーバ 1 は、ユーザの質問文に対して複数の回答文の候補が存在する場合、収集したコメント文章を用いて、ユーザの意図を反映した回答文を明示する(に絞り込む)。

【 0 0 2 8 】

図 2 は、本発明における質問回答サーバの機能構成図である。

【 0 0 2 9 】

図 2 によれば、質問回答サーバ 1 は、通信インタフェース部 1 0 と、回答文蓄積部 1 0 1 と、回答文取得部 1 1 1 と、コメント文章蓄積部 1 0 2 と、コメント文章収集部 1 1 2 とを有する。

【 0 0 3 0 】

回答文蓄積部 1 0 1 は、多数の回答文を蓄積する。回答文取得部 1 1 1 が、これら回答文を、ネットワークを介して回答文蓄積サーバ 2 から受信し、回答文蓄積部 1 0 1 へ蓄積するものであってもよい。

【 0 0 3 1 】

コメント文章蓄積部 1 0 2 は、不特定多数の第三者によって投稿された多数のコメント文章を蓄積する。コメント文章収集部 1 1 2 が、これらコメント文章を、ネットワークを介してブログサーバ 3 から受信し、コメント文章蓄積部 1 0 2 へ蓄積するものであってもよい。

【 0 0 3 2 】

「コメント文章」とは、例えば twitter (登録商標) で発信された、日本語の「つぶやき」(最大文字数: 140 文字)のようなものである。コメント文章は、例えば、ユーザ id(from_user_id)、つぶやき ID(id_str)、発信時間(created_at)、つぶやき(texts)を含む。ここで、コメント文章収集部 1 1 2 は、予め指定した複数のキーワードを含むコメント文章のみを収集することもできる。

【 0 0 3 3 】

また、図 2 によれば、質問回答サーバ 1 は、質問文入力部 1 2 1 と、質問キーワード抽出部 1 2 2 と、コメント文章検索部 1 2 3 と、トピック分類部 1 2 4 と、回答文検出部 1 2 5 と、代表述語項構造抽出部 1 2 6 と、差分述語項構造抽出部 1 2 7 と、差分述語項構造選択部 1 2 8 と、回答文出力部 1 2 9 とを有する。これら機能構成部は、サーバに搭載されたコンピュータを機能させるプログラムを実行することによって実現される。

【 0 0 3 4 】

[質問文入力部 1 2 1]

質問文入力部 1 2 1 は、質問者の端末 4 から、ネットワークを介して質問文を受信する。例えばユーザの質問文は、以下のようなものである。

Q「携帯電話機の紛失」

その質問文は、質問キーワード抽出部 1 2 2 へ出力される。

【 0 0 3 5 】

図 3 は、質問キーワード抽出部及びコメント文章検索部の処理を表す説明図である。

【 0 0 3 6 】

[質問キーワード抽出部 1 2 2]

質問キーワード抽出部 1 2 2 は、質問文に含まれる複数の質問キーワードを抽出する。ここで、質問キーワード抽出部 1 2 2 は、質問文から形態素解析によってキーワードを抽出すると共に、TF-IDF (Term Frequency - Inverse Document Frequency: 単語の出現頻度 - 逆出現頻度) によって特徴的な単語を、質問キーワードとして抽出する。

【 0 0 3 7 】

質問キーワード抽出部 1 2 2 は、最初に、質問文から形態素解析によって単語を抽出する。「形態素解析」とは、文章を、意味のある単語に区切り、辞書を利用して品詞や内容を判別する技術をいう。「形態素」とは、文章の要素のうち、意味を持つ最小の単位を意味する。形態素解析のように単語単位で検索することなく、文字単位で分解し、後続の N-1 文字を含めた状態で出現頻度を求める「N-gram」によって解析するものであってもよい。

【 0 0 3 8 】

次に、TF-IDF によって特徴的なキーワードを、質問キーワードとして抽出する。TF-IDF とは、各単語に重みを付けて、クエリから文章をベクトル空間で表し、文章とクエリの類似度でランク付けをする技術である。ランク付けられた値が高いほど、重要キーワードと認識される。

【 0 0 3 9 】

図 3 の例によれば、以下のように抽出される。

質問文 「携帯電話機の紛失」

質問キーワード「携帯電話機」「紛失」

【 0 0 4 0 】

[コメント文章検索部 1 2 3]

コメント文章検索部 1 2 3 は、コメント文章蓄積部 1 0 2 を用いて、質問キーワードを含むコメント文章を検索する。具体的には、質問キーワードをクエリとして、各コメント文章から TF (Term Frequency) 値や DF (Document Frequency) 値を抽出し、これら値が所定閾値以上となる複数のコメント文章を検索する。TF 値は、文章における検索語の出現頻度をいい、DF 値は、索引語が現れる相対文章頻度をいう。コメント文章検索部 1 2 3 は、ソーシャルメディア検索機能であって、投稿された大量のつぶやきの中から、質問キーワードに関するつぶやきのみを検索するようなものである。

【 0 0 4 1 】

図 3 によれば、例えば 4 つのコメント文章が検索されている。これらコメント文章には、少なくとも「携帯電話機」又は「紛失」が含まれている。

【 0 0 4 2 】

図 4 は、トピック分類部の処理を表す説明図である。

【 0 0 4 3 】

[トピック分類部 1 2 4]

トピック分類部 1 2 4 は、検索された複数のコメント文章を、述語項構造解析によって、述語項構造の分布から複数個のトピックグループに分類する。トピック分類部 1 2 4 は、当該コメント文章を、分類された各トピックグループに属する確からしさ(トピック比率)を算出する LDA (Latent Dirichlet Allocation) アルゴリズムを用いて、いずれか

10

20

30

40

50

1つのトピックグループに分類する。特に、トピック分類部124のLDAは、キーワードによる分類でなく、述語項構造による分類である。

【0044】

LDAは、単語文書行列を次元圧縮する技術(LSI(latent Semantic Indexin))に対して、単語の特徴ベクトルに揺らぎに基づく確率的な枠組みを導入したものである(例えば非特許文献2参照)。その圧縮した次元の集合をトピックという。

【0045】

また、「述語項構造」とは、文章中の述語に対して「項」となる名詞句等を当てたものである。述語項構造を用いることによって、文章の意味の骨格を把握することができる。述語項構造解析として、例えばフリーソフトであるSyncha等の述語項構造解析器を用いることができる。

10

【0046】

述語項構造は、「述語」に対する「目的語」とその格とから構成される。図4によれば、例えば「携帯を探す」の述語項構造は、述語「探す」に対して目的語「携帯」及び格「ヲ」からなる。また、例えば「サービスに申し込む」の述語項構造は、述語「申し込む」に対して目的語「サービス」及び格「ニ」からなる。

【0047】

トピック分類部124は、以下のステップで処理を実行する。

(S41) 質問キーワードに関する多数のコメント文章から、述語項構造毎の出現頻度(出現回数)をLDA処理へ入力する。そして、コメント文章毎に、各述語項構造の出現頻度を計数する。

20

(S42) 次に、本件でのLDA処理では、トピック毎の述語項構造分布や、コメント文章(ネット側意見)毎のトピック比率を取得する。このトピック比率によって、コメント文章が属するトピックグループに分類する。そして、トピックグループ毎に、全てのコメント文章に含まれる各述語項構造の出現頻度を計数する。

(S43) 次に、コメント文章毎に、各トピックグループに属する述語項構造を計数する。そして、コメント文章を計数値の高いトピックグループに分類する。

【0048】

図5は、回答文検出部の処理を表す説明図である。

【0049】

30

[回答文検出部125]

回答文検出部125は、各トピックグループに含まれるコメント文章群と、各回答文に含まれる文章との間の類似度を算出し、各トピックグループに類似度が所定閾値以上となる回答文を対応付ける。

【0050】

類似度の算出方法は、例えば以下のようにする。

(S51) 回答文検出部125は、各トピックグループに含まれるコメント文章群から述語項構造解析によって述語項構造を抽出すると共に、トピックグループにおける第1の特徴ベクトルを算出する。

各トピックグループ: $C_i (i=1, 2, \dots)$

40

トピックグループiに含まれるコメント文章: $T_{ij} (j=1, 2, \dots)$

(S52) 回答文蓄積部101に蓄積された各回答文から述語項構造解析によって述語項構造を抽出すると共に、当該回答文における第2の特徴ベクトルとを算出する。

回答文: $A_k (k=1, 2, \dots)$

(S53) トピックグループの第1のベクトルと、回答文の第2のベクトルとの間のコサイン距離に基づいて類似度を算出する。具体的には、各コメント文章 T_{i1}, T_{i2}, \dots を含むトピックグループ C_i と、回答文 A_j との類似度 $\text{Dist}(C_i, A_j)$ を算出する。

$\text{Dist}(C_i, A_j) = \cos \text{in 距離 } D(T_{i1}, A_j), D(T_{i2}, A_j), \dots \text{ の平均値}$
 $= \arg \max (\text{Dist}(C_i, A_j))$

【0051】

50

図 5 によれば、回答文蓄積部 101 には、多数の回答文が蓄積されている。

回答文 1 「・・・」

回答文 2 「携帯探せて安心サービスの申込方法」

述語項構造 = 述語「探す」、目的語「携帯」、ヲ格

述語項構造 = 述語「申し込む」、目的語「サービス」、二格

回答文 3 「・・・」

回答文 4 「・・・」

回答文 5 「携帯探せて安心サービスの利用方法」

述語項構造 = 述語「探す」、目的語「携帯」、ヲ格

述語項構造 = 述語「利用する」、目的語「サービス」、ヲ格

10

回答文 6 「・・・」

【0052】

図 5 によれば、トピックグループ 1 と回答文 2 との類似度が、所定閾値 よりも高い場合、両者は類似していると判定されている。また、トピックグループ 2 と回答文 5 との類似度が、所定閾値 よりも高い場合、両者は類似していると判定されている。これによって、トピックグループ C1, C2, ... 毎に、0 個以上の回答文が割り当てられる。

【0053】

図 6 は、代表述語項構造抽出部、差分述語項構造抽出部、差分述語項構造選択部及び回答文出力部の処理を表す説明図である。

【0054】

20

[代表述語項構造抽出部 126]

代表述語項構造抽出部 126 は、各トピックグループについて、対応付けられた回答文に含まれる述語項構造の中で、当該トピックグループを特徴付ける代表述語項構造を抽出する。

【0055】

図 6 によれば、トピックグループ 1 に対応する回答文 2 からは、以下の表のような述語項構造が抽出される。

回答文 2 「携帯探せて安心サービスの申し込み方法は以下ようになります・・・」

述語項構造 = 述語「探す」、目的語「携帯」、ヲ格

述語項構造 = 述語「申し込む」、目的語「サービス」、二格

30

回答文 5 「携帯探せて安心サービスの利用方法は以下ようになります・・・」

述語項構造 = 述語「探す」、目的語「携帯」、ヲ格

述語項構造 = 述語「利用する」、目的語「サービス」、ヲ格

【0056】

ここで、述語項構造抽出部 126 は、各トピックグループの代表述語項構造を、赤池情報量基準に応じて優先順に並べることも好ましい。トピックグループ C1, C2, ... に割り当てられた回答文のいずれかに出現する述語項構造を、s1, s2, ... とする。ここでは、述語項構造 E(i) が、トピックグループ Cj の判別に役立つかどうかの指標を与える。

【0057】

以下では、述語項構造 s が、トピックグループ C の判別に役立つかどうかの指標 E (s, C) の算出方法を表す。

40

【0058】

(S1) トピックグループに含まれるコメント文章 (つぶやき) の集合 U から、以下の 4 種類の頻度を得る。

n11 = トピックグループ C に類似し、述語項構造 s が出現するコメント文章の数

n12 = トピックグループ C 以外に類似し、述語項構造 s が出現するコメント文章の数

n21 = トピックグループ C に類似し、述語項構造 s が出現しないコメント文章の数

n22 = トピックグループ C 以外に類似し、述語項構造 s が出現しないコメント文章の数

【0059】

(S2) 次に、n11, n12, n21, n22 に対して、赤池情報量規準 (AIC : Akaike's Infor

50

mation Criterion)を用いて、独立モデルに対する値 $MLL_IM(s,C)$ 及び従属モデルに対する値 $MLL_DM(s,C)$ を算出する。これは、述語項構造とトピックグループとの組毎の不当割合を算出する。

$$\begin{aligned} MLL_IM(s,C) = & (n_{11}+n_{12}) \log(n_{11}+n_{12}) \\ & + (n_{11}+n_{21}) \log(n_{11}+n_{21}) \\ & + (n_{21}+n_{22}) \log(n_{21}+n_{22}) \\ & + (n_{12}+n_{22}) \log(n_{12}+n_{22}) - 2N \log N \end{aligned}$$

$$MLL_DM(s,C) = n_{11} \log n_{11} + n_{12} \log n_{12} + n_{21} \log n_{21} + n_{22} \log n_{22} - N \log N$$

但し、 $N = n_{11} + n_{12} + n_{21} + n_{22}$

【 0 0 6 0 】

10

(S 3) 前述の $MLL_IM(s,C)$ 及び $MLL_DM(s,C)$ から、以下の $E(s,C)$ を算出する。

$$AIC_IM(s,C) = -2 \times MLL_IM(s,C) + 2 \times 2$$

$$AIC_DM(s,C) = -2 \times MLL_DM(s,C) + 2 \times 3$$

$$E(s,C) = AIC_IM(s,C) - AIC_DM(s,C)$$

【 0 0 6 1 】

前述で算出された $E(s,C)$ は、述語項構造 s がトピックグループ C に偏って出現する不当割合を表す。 $E(s,C)$ は、赤池情報量基準に従って、トピックグループ C の判別に役立つ述語項構造ほど、 $E(s,C)$ の値が高くなる。本発明によれば、各トピックグループ C_i に対し、 $E(s,C)$ の値が大きい順に、 m 個の述語項構造 $C_{i,1}$ 、 $C_{i,2}$ 、 $C_{i,3}$ 、 \dots 、 $C_{i,m}$ を抽出し、トピックグループ C_i の代表述語項構造とする。

20

【 0 0 6 2 】

[差分述語項構造抽出部 1 2 7]

差分述語項構造抽出部 1 2 7 は、各トピックグループについて、当該トピックグループのみに出現する代表述語項構造を、差分述語項構造として抽出する。回答文 2 及び 5 について、[述語項構造 = 述語「探す」、目的語「携帯」、ヲ格] は共通する。そこで、図 6 によれば、以下の差分述語項構造が抽出される。

回答文 2 「携帯探せて安心サービスの申し込み方法は以下ようになります・・・」

述語項構造 = 述語「申し込む」、目的語「サービス」、二格

回答文 5 「携帯探せて安心サービスの利用方法は以下ようになります・・・」

述語項構造 = 述語「利用する」、目的語「サービス」、ヲ格

30

【 0 0 6 3 】

[差分述語項構造選択部 1 2 8]

差分述語項構造選択部 1 2 8 は、複数の差分述語項構造に基づく文章を、ユーザインタフェースを介してユーザに明示する。

【 0 0 6 4 】

差分述語項構造からの日本語文章を生成するために、例えば以下のようなルールが設定される。

(ルール 1) ヲ格のみからなる述語項構造 S の場合

$$\rightarrow W(S, \text{ヲ格}) + \text{「を」} + W(S, \text{述語})$$

(ルール 2) ヲ格とデ格からなる述語項構造 S の場合

$$\rightarrow W(S, \text{デ格}) + \text{「で」} + W(S, \text{ヲ格}) + \text{「を」} + S(\text{述語})$$

40

(ルール 3) 差分述語項構造 S_1 のヲ格と、共通述語項構造 B のデ格とが一致する場合、

$$\rightarrow A \text{ から生成した日本語のヲ格の前方修飾語として、} B \text{ の日本語を埋め込む}$$

$W(S, \text{ヲ格})$ は、述語項構造 S のヲ格の単語を表す。

$W(S, \text{デ格})$ は、述語項構造 S のデ格の単語を表す。

尚、このようなルールに基づく日本語の生成については、機械翻訳システムの技術が適用できる(例えば非特許文献 3 参照)。

【 0 0 6 5 】

例えば、ルール 1 ~ 3 を用いて、図 6 によれば、トピックグループ毎に、以下の 2 つの日本語文章が生成される。

50

「携帯を探すサービスを申し込む」

「携帯を探すサービスを利用する」

【0066】

これに対し、端末4は、ユーザ操作に応じていずれか1つの差分述語項構造の文章を選択させる。ユーザから見ると、例えば、質問文をキーボードで入力した後、トピックグループ毎の差分述語項構造に基づいた自然な日本語文章がディスプレイに表示される。そして、ユーザは、いずれかの文章を選択することができる。ここで、図6によれば、ユーザは、「携帯を探すサービスを利用する」を選択している。ユーザに選択された文章の差分述語項構造は、回答文出力部129へ出力される。

【0067】

10

[回答文出力部129]

回答文出力部129は、回答文検出部125によって検出された回答文を、対応する1つ以上の差分述語項構造の文章と共に明示する。本発明によれば、ユーザの質問に曖昧性があり、コメント文章群が複数のトピックグループに分類され、各トピックグループに対応付けられた回答文を得ることができる。ここで、この得られた回答文の数が少ない場合、差分述語項構造に基づく文章は、提示される回答文の傾向をユーザが認識するために有益な情報となる。

【0068】

また、回答文の数が多い場合、ユーザとインタラクション（やりとり）をすることによって、回答文を絞り込むことが好ましい。そこで、回答文出力部129は、選択された差分述語項構造に対応する回答文を、ユーザインタフェースを介して明示する。例えば、その回答文を、ユーザが視認するディスプレイに表示する。図6によれば、「携帯探せて安心サービスの利用方法」の回答文が、ユーザへ表示される。これによって、ユーザは、質問文に対する回答文を認識することができる。

20

【0069】

図7は、本発明におけるシーケンス図である。

【0070】

(S71) 質問者が操作する端末4から、質問回答サーバ1へ、ユーザの質問文が送信される(図2の質問文入力部121参照)。

(S72) 質問回答サーバ1は、質問文に含まれる複数の質問キーワードを抽出する(図2の質問キーワード抽出部122参照)。

30

(S73) 質問回答サーバ1は、コメント文章蓄積部102を用いて、質問キーワードを含むコメント文章を検索する(図2のコメント文章検索部123参照)。

(S74) 質問回答サーバ1は、検索された複数のコメント文章を、述語項構造解析によって、述語項構造の分布から複数個のトピックグループに分類する(図2のトピック分類部124参照)。

(S75) 質問回答サーバ1は、各トピックグループに含まれるコメント文章群と、各回答文との間の類似度を算出し、各トピックグループに類似度が所定閾値以上となる回答文を対応付ける(図2の回答文検出部125参照)。

(S76) 質問回答サーバ1は、各トピックグループについて、対応付けられた回答文に含まれる述語項構造の中で、当該トピックグループを特徴付ける代表述語項構造を抽出する(図2の代表述語項構造抽出部126参照)。

40

(S77) 質問回答サーバ1は、各トピックグループについて、当該トピックグループのみに出現する代表述語項構造を、差分述語項構造として抽出する(図2の差分述語項構造抽出部127参照)。

(S78) 質問回答サーバ1は、複数の差分述語項構造に基づく文章を、ユーザ操作の端末4へ送信する(図2の差分述語項構造選択部128参照)。そして、端末4では、ユーザ操作に応じていずれか1つの文章が選択させる。選択された文章の差分述語項構造は、端末4から質問回答サーバ1へ送信される。

(S79) 質問回答サーバ1は、選択された差分述語項構造に対応する回答文を、ユーザ

50

の端末 4 へ送信する（図 2 の回答文出力部 1 2 9 参照）。

【0071】

前述したように本発明の質問回答サーバによれば、例えばtwitterのような大量のコメント文章から、質問文の意図を表す代表的な述語項構造を抽出し、質問文を補完することによって、回答文を高精度に検索することができる。具体的には、最初に、質問文に含まれるキーワードを抽出してソーシャルメディアを検索し、大量の検索結果を複数のトピックグループ（トピック毎に1つの検索意図に対応）に高速に分類し、各トピックに類似する回答文を回答文蓄積部から検索する。次に、各トピックグループに特有の単語（差分述語項構造）を自動的に抽出してユーザに提示し、ユーザの選択結果に従った回答文に絞り込んで、ユーザとの対話形式を繰り返し実行することができる。

10

【0072】

以上、詳細に説明したように、本発明の質問回答プログラム、サーバ及び方法によれば、ユーザの質問文に対して複数の回答文の候補が存在する場合、ユーザの意図を反映した回答文を明示する（に絞り込む）ことができる。

【0073】

最後に、本発明が、キーワード検索ではなく、述語項構造検索を用いた効果について詳述する。

【0074】

一般に、例えば、「携帯電話が紛失したらどうしよう」というユーザからの短い質問の場合、「携帯電話紛失に備えたサービス申し込みの要望」なのか、又は、「端末の紛失への対応法に関する問い合わせ」なのかといった曖昧性が存在する。これに対し、コンテキストに依存した曖昧性を検出し、対話形式で回答候補を絞り込みながらTIPS等を返答する技術が提案されている（本願と同一出願人及び同一発明者によって出願された平成24年1月12日付け特許出願、以下「先の出願に係る発明」と称す）。この技術によれば、コンテキストの曖昧性を高速かつ適切に検出し、「安心、申し込み」や「端末、発見」といったキーワードを対話の選択肢として提示することができる。

20

【0075】

しかしながら、先の出願に係る発明によれば、第1の課題として、コンテキストを絞り込む際に「携帯紛失、サービス、申し込み」又は「携帯紛失、サービス、利用」といったキーワードが提示されるだけであって、ユーザにとっては、コンテキストの差異を理解しづらい。

30

【0076】

また、第2の課題として、コンテキストの絞り込みが終了しても、非特許文献3に記載された技術のようなキーワードによる検索によれば、検索条件としての情報が不足することがある。例えば、「携帯＋発見」といったキーワードで検索する場合、宝探しゲームのように携帯電話機で何かを発見するサービスや、携帯電話を発見するサービスを発見するサービスの情報の両方がキーワード検索結果に現れる。このため、回答精度を現状以上に向上させることが難しい。

【0077】

第1の課題に対して、本発明によれば、差分キーワードの単純な提示ではなく、動詞を中心に主語や目的語等の関係をリンクで表す「述語項構造」と呼ばれるデータから自然な応対の文章を生成し、それを利用者に提示する手法を用いることで理解度が深まることが期待できる。例えば、「携帯＋紛失」のユーザクエリに対して従来システムが「安心＋申し込み」もしくは、「発見」といったキーワードを提示していたのに対し、「端末を紛失した際に安心できる申し込みに関する情報」もしくは「携帯電話を紛失した際に端末を発見すること」といった自然な文章をユーザに提示することでユーザの利便性を向上させる手段を提供する。

40

【0078】

第2の課題に対して、本発明によれば、述語項構造を検索パラメータとして知識源のテキストを検索することにより、コンテキスト絞り込み後の検索精度を大きく向上させる。

50

例えば、「端末を発見する」と「端末で発見する」の意味を区別できる述語項構造を用いるので、従来型のキーワード検索で行われていた過剰検出が減る。

【 0 0 7 9 】

また、述語項構造を使用して文書の類似性を判定する場合、抽出した述語項構造の一致度合いを判定する必要がある、従来のキーワードを利用した場合より多くの計算時間が必要となることも問題となる。この問題に対しては、述語項構造の中で格と呼ばれるデータスロットに注目し、使用頻度が高いスロットの組み合わせを事例から事前に学習しておき、使用頻度の高い組み合わせに対してはハッシュ関数を使って高速に検索できるようにする。

【 0 0 8 0 】

前述した本発明の種々の実施形態について、本発明の技術思想及び見地の範囲の種々の変更、修正及び省略は、当業者によれば容易に行うことができる。前述の説明はあくまで例であって、何ら制約しようとするものではない。本発明は、特許請求の範囲及びその均等物として限定するものにのみ制約される。

【 符号の説明 】

【 0 0 8 1 】

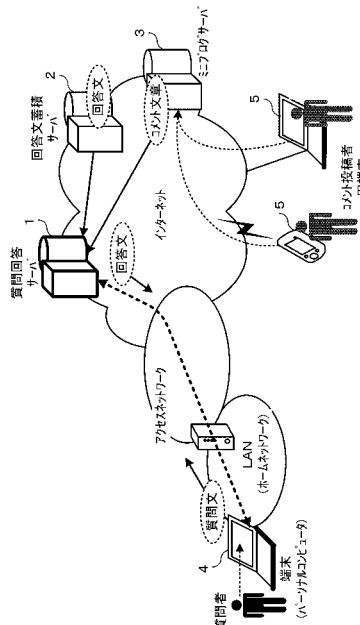
- 1 質問回答サーバ
- 1 0 通信インタフェース部
- 1 0 1 回答文蓄積部
- 1 0 2 コメント文章蓄積部
- 1 1 1 回答文取得部
- 1 1 2 コメント文章収集部
- 1 2 1 質問文入力部
- 1 2 2 質問キーワード抽出部
- 1 2 3 コメント文章検索部
- 1 2 4 トピック分類部
- 1 2 5 回答文検出部
- 1 2 6 代表述語項構造抽出部
- 1 2 7 差分述語項構造抽出部
- 1 2 8 差分述語項構造選択部
- 1 2 9 回答文出力部
- 2 回答文蓄積サーバ
- 3 ブログサーバ
- 4 端末
- 5 コメント投稿者用の汎用端末

10

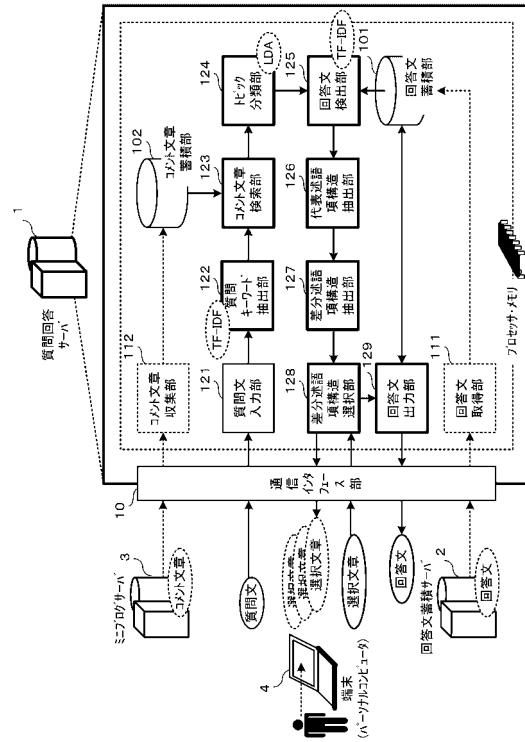
20

30

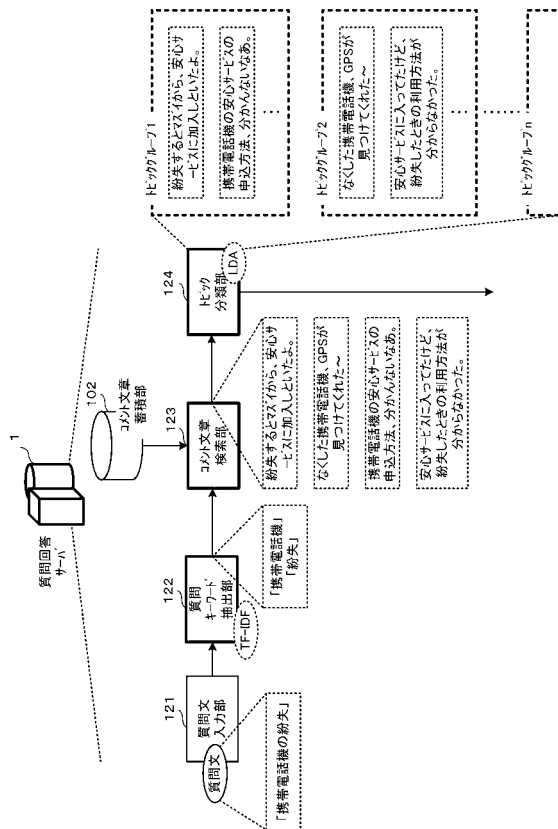
【 図 1 】



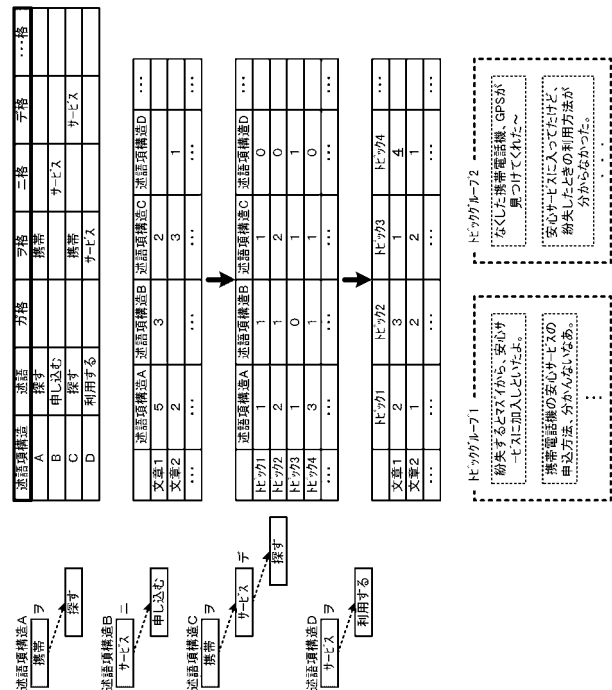
【 図 2 】



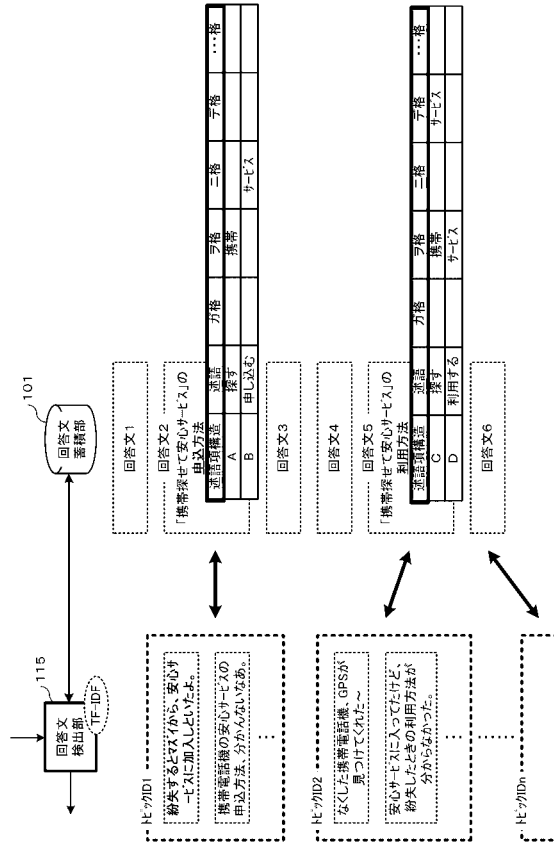
【 図 3 】



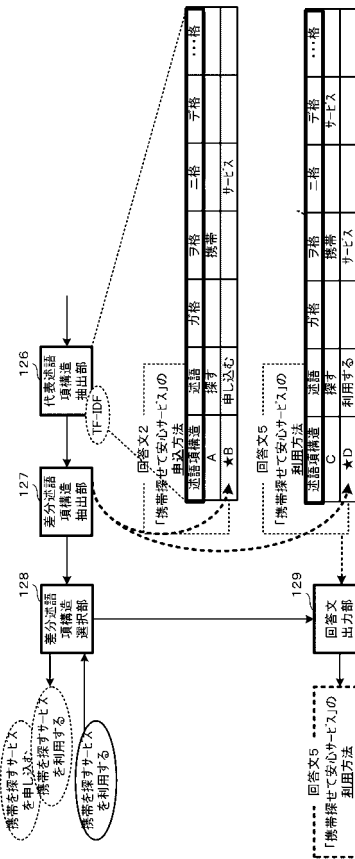
【 図 4 】



【 図 5 】



【 図 6 】



【 図 7 】

