

(19) 日本国特許庁(JP)

(12) 特 許 公 報(B2)

(11) 特許番号

特許第4904496号
(P4904496)

(45) 発行日 平成24年3月28日(2012.3.28)

(24) 登録日 平成24年1月20日(2012.1.20)

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

請求項の数 8 (全 27 頁)

| | |
|---|--|
| <p>(21) 出願番号 特願2006-304301 (P2006-304301) (22) 出願日 平成18年11月9日(2006.11.9) (65) 公開番号 特開2008-123111 (P2008-123111A) (43) 公開日 平成20年5月29日(2008.5.29) 審査請求日 平成21年9月25日(2009.9.25)</p> <p>特許法第30条第1項適用 2006年5月19日 社 団法人 情報処理学会発行の「情報処理学会研究報告 2006-NL-173」に発表</p> | <p>(73) 特許権者 504174135 国立大学法人九州工業大学 福岡県北九州市戸畑区仙水町1番1号</p> <p>(74) 代理人 100099634 弁理士 平井 安雄</p> <p>(72) 発明者 野村 浩郷 福岡県飯塚市大字川津680-4 九州工 業大学内</p> <p>審査官 鈴木 和樹</p> |
|---|--|

最終頁に続く

(54) 【発明の名称】 文書類似性導出装置及びそれを用いた回答支援システム

(57) 【特許請求の範囲】

【請求項1】

文からなる文書の文を形態素解析する手段と、

形態素解析された文書から、当該文書に出現する索引語のTF / IDFによる重みを要素としたTF / IDFベクトルを求める手段と、

形態素解析された文書から、当該文書に出現する体言について当該体言が出現する文中で共起した用言の頻度を要素とした共起ベクトルを求める手段とを含み、

第1の文書のTF / IDF文書ベクトル及び共起ベクトルを求め、

第2の文書のTF / IDF文書ベクトル及び共起ベクトルを求め、

求めた第1の文書のTF / IDF文書ベクトル及び共起ベクトルと第2の文書のTF / IDF文書ベクトル及び共起ベクトルから第1の文書と第2の文書の類似性を求める文書類似性導出装置。

10

【請求項2】

文からなる文書の文章を形態素解析する手段と、

形態素解析された文書から、当該文書に出現する索引語のTF / IDFによる重みを要素としたTF / IDFベクトルを求める手段と、

形態素解析された文書から、当該文書に出現する索引語が出現する文章の文タイプを決定し、それぞれの文タイプの頻度を要素とした文タイプベクトルを求める手段とを含み、

第1の文書のTF / IDF文書ベクトル及び文タイプベクトルを求め、

第2の文書のTF / IDF文書ベクトル及び文タイプベクトルを求め、

20

求めた第1の文書のTF/IDF文書ベクトル及び文タイプベクトルと第2の文書のTF/IDF文書ベクトル及び文タイプベクトルから第1の文書と第2の文書の類似性を求める文書類似性導出装置。

【請求項3】

文からなる文書の文を形態素解析する手段と、

形態素解析された文書から、当該文書に出現する索引語のTF/IDFによる重みを要素としたTF/IDFベクトルを求める手段と、

形態素解析された文書から、当該文書に出現する体言について当該体言が出現する文中で共起した用言の頻度を要素とした共起ベクトルを求める手段と、

形態素解析された文書から、当該文書に出現する索引語が出現する文章の文タイプを決定し、それぞれの文タイプの頻度を要素とした文タイプベクトルを求める手段とを含み、

第1の文書のTF/IDF文書ベクトル、共起ベクトル及び文タイプベクトルを求め、

第2の文書のTF/IDF文書ベクトル、共起ベクトル及び文タイプベクトルを求め、

求めた第1の文書のTF/IDF文書ベクトル、共起ベクトル及び文タイプベクトルと第2の文書のTF/IDF文書ベクトル、共起ベクトル及び文タイプベクトルから第1の文書と第2の文書の類似性を求める文書類似性導出装置。

【請求項4】

前記請求項1ないし3のいずれかに記載の文書類似性導出装置の各手段を含み、

TF-IDFベクトルと共起ベクトル及び/又は文タイプベクトルである文書ベクトルを第1の文書について求め、

複数文書からなる第2の文書群の各文書の文書ベクトルを求め、

求めた第2の文書群の各文書の文書ベクトルから平均文書ベクトルを求め、

求めた第2の文書群の平均文書ベクトルと第1の文書の文書ベクトルから第1の文書と第2の文書群の類似性を求める文書-文書群類似性導出装置。

【請求項5】

前記請求項1ないし3のいずれかに記載の文書類似性導出装置の各手段を含み、

TF-IDFベクトルと共起ベクトル及び/又は文タイプベクトルである文書ベクトルを比較対象となる比較対象文書について求め、

第nの文書の索引TF-IDFの文書ベクトル、共起ベクトル及び文タイプベクトルを求め、

比較対象文書の文書ベクトルと第nの文書の文書ベクトルから比較対象文書と第nの文書の類似性を求め、

nは1ないしNまでであり、各第nの文書と比較対象文書の類似性の中から類似性の高い第nの文書を特定する高類似性文書特定装置。

【請求項6】

前記請求項1ないし3のいずれかに記載の文書類似性導出装置の各手段を含み、

TF-IDFベクトルと共起ベクトル及び/又は文タイプベクトルである文書ベクトルを比較対象文書について求め、

複数文書からなる第nの文書群の各文書の文書ベクトルを求め、

求めた第nの文書群の各文書の文書ベクトルから平均文書ベクトルを求め、

求めた第nの文書群の平均文書ベクトルと第1の文書の文書ベクトルから第1の文書と第nの文書群の類似性を求め、

nは1ないしNまでであり、各第nの文書群と比較対象文書の類似性の中から類似性の高い第nの文書群を特定する高類似性文書群特定装置。

【請求項7】

文からなる文書の文を形態素解析する手段と、

形態素解析された文書から、当該文書に出現する索引語のTF/IDFによる重みを要素としたTF/IDFベクトルを求める手段と、

形態素解析された文書から、当該文書に出現する体言について当該体言が出現する文中で共起した用言の頻度を要素とした共起ベクトルを求める手段と、

10

20

30

40

50

形態素解析された文書から、当該文書に出現する索引語が出現する文章の文タイプを決定し、それぞれの文タイプの頻度を要素とした文タイプベクトルを求める手段としてコンピュータを機能させ、

第1の文書のTF/IDF文書ベクトル、共起ベクトル及び文タイプベクトルを求め、
第2の文書のTF/IDF文書ベクトル、共起ベクトル及び文タイプベクトルを求め、
求めた第1の文書のTF/IDF文書ベクトル、共起ベクトル及び文タイプベクトルと第2の文書のTF/IDF文書ベクトル、共起ベクトル及び文タイプベクトルから第1の文書と第2の文書の類似性をコンピュータに求めさせる文書類似性導出プログラム。

【請求項8】

文からなる文書の文を形態素解析するステップと、
形態素解析された文書から、当該文書に出現する索引語のTF/IDFによる重みを要素としたTF/IDFベクトルを求めるステップと、

形態素解析された文書から、当該文書に出現する体言について当該体言が出現する文中で共起した用言の頻度を要素とした共起ベクトルを求めるステップと、

形態素解析された文書から、当該文書に出現する索引語が出現する文章の文タイプを決定し、それぞれの文タイプの頻度を要素とした文タイプベクトルを求めるステップとを含み、

第1の文書のTF/IDF文書ベクトル、共起ベクトル及び文タイプベクトルを求めるステップと、

第2の文書のTF/IDF文書ベクトル、共起ベクトル及び文タイプベクトルを求めるステップと、

求めた第1の文書のTF/IDF文書ベクトル、共起ベクトル及び文タイプベクトルと第2の文書のTF/IDF文書ベクトル、共起ベクトル及び文タイプベクトルから第1の文書と第2の文書の類似性を求めるステップとをさらに含む文書類似性導出方法。

【発明の詳細な説明】

【技術分野】

【0001】

本発明は、文書の類似性を求める文書類似性導出装置に関する。

【背景技術】

【0002】

近年インターネットやパソコンの普及により、アフターサービスの一環としてパソコン技術サポートの必要性が高まっている。多くのパソコン技術サポートセンターでは、主に電話で技術サポートを行う従来型のコールセンターに加えて、インターネット経由でE-mailでの問い合わせを受け付けるメールコールセンターがたくさん設置されてきている。

【0003】

メールコールセンターで行われている技術サポートは、すべて無料サポートである。質問メールは、夕刻から深夜にかけて多く送付されてくる。問い合わせメールの受信から回答の発信までは所定時間内（例えば24時間以内）に完了することが求められている。このような制約があるため、企業にとって、正確かつ迅速なサポートを行うには人件費などのコストが膨大なものになりつつある。そこで、メールコールセンターの自動化が強く求められている。

そこで、問い合わせに適した回答作成を支援する回答支援装置が、特開2001-273308号公報に開示されている。

【0004】

この特開2001-273308号公報の回答支援装置は、予め想定された問い合わせの内容とこの問い合わせに対する回答作成を支援する支援情報とが対応づけられて格納される支援情報記憶手段と、予め想定された問い合わせの内容とこの問い合わせに対する回

10

20

30

40

50

答作成者の回答作成者情報とが対応づけられて格納される回答作成者情報記憶手段と、入力される問い合わせの内容により前記支援情報記憶手段を検索して得られた支援情報および当該問い合わせの内容を、当該内容により前記回答作成者情報記憶手段を検索して得られた回答作成者情報に対応する回答作成者に送付する情報送付手段とを有するものである。

【特許文献1】特開2001-273308号公報

【発明の開示】

【発明が解決しようとする課題】

【0005】

前記背景技術の回答支援装置によれば、消費者からの問い合わせに対して適切な回答作成者に問い合わせが転送されると共に、回答作成者には問い合わせに係る支援情報を得ることができるため迅速に回答することができる。

10

【0006】

しかしながら、この背景技術の回答支援装置は、具体的には、消費者がリストボックスやチェックボックスなどの選択形式の問い合わせに対して対応する回答作成者及び支援情報を特定するものであり、消費者が自ら作成した問い合わせ文章に対応することができないという課題を有する。なお、支援情報とは、具体的には、消費者への問い合わせに係るマニュアル、仕様書のことである。

【0007】

本発明は前記課題を解決するためになされたものであり、文章による問い合わせに対してその種類を特定し、適切な回答を支援する回答支援システムを提供することを目的とする。また、この回答支援システムで用いる文書間の類似性を求める新たな手法を提供することも目的とする。

20

【課題を解決するための手段】

【0008】

消費者などの質問者からの質問に対して、システムが少数の回答候補を提示して、それらから回答者が最適なものを選択し、最終査読して回答を仕上げるような支援システムを作成し、使用するのが現実的である。このようなシステムでは、最終査読以外は自動的に処理されることになり、コストの大幅な削減が実現できる。発明者は鋭意努力によりシステムACCESS(Automated Call-Center Service System)を作成した。

30

【0009】

実際のメールコールセンターでの約三年間の実務により収集した三万件以上の最終査読済みの「質問応答」データの中から、約一万件を使って「質問応答データベース」を構築した。

【0010】

実際の質問には、同一内容のものや似た内容のものが多い。したがって、ユーザからの質問が来る度に逐一回答を作成するのは無駄である。そのため、既に回答した質問応答データから「質問応答データベース」を構築し、それを再利用できるようにすると、大幅なコストダウンができる。

40

【0011】

「ユーザからの質問」には、「質問」とは思われないものも含まれている。このような「質問」に対しては、最終査読者が「査読」するのではなく、異なる観点からの対応が必要になるものもある。

【0012】

ユーザからの質問は自由記述であるので、質問文にはミスタイプ、かな漢字変換の誤り、文法的不完全さなどが多く見られる。このような現象を前提として処理しなければならない。しかし、最終査読済みの質問応答データではそれらは修正・訂正されているので、質問応答データベースを検索して、うまく活用することには大きな利点がある。

【0013】

50

質問応答データベースは、二層にカテゴリ化した構造としている。それぞれのカテゴリにはそれぞれを特徴づけるタグを付与している。最下層のカテゴリには、実際の質問応答データが格納されている。二層にカテゴリ化した理由は、最下層でうまくマッチするものがないとき、いわゆるシソーラスにおける上位概念を利用するという考えに似ている。なお、ここで二層を示したが三層以上であってもよい。

【0014】

質問応答データベースの再利用に関しては、パソコンユーザから送られてきた問い合わせメールから、その質問がどの質問カテゴリに属するのかを統計的処理などにより推定することにより行う。推定された質問カテゴリごとのカテゴリ回答文を用いて作成した回答候補を最終査読者に提示することで、回答作成を効率化することになる。

10

【0015】

質問カテゴリ推定の精度を検証するために、システムを作成し、評価実験を行った。145個の質問カテゴリを用いて実験したところ、86%の割合で正解の質問カテゴリを上位3位以内に推定した。蓄積されている三万件以上の「質問応答」データの中から、約一万件を使って「質問応答データベース」を構築してみたため、残りの約二万件の「質問応答」データを「質問応答データベース」に加えることにより、正解の質問カテゴリを発見する精度は大幅に改善できる。

【0016】

(1) 文書間の類似性(図1、図2参照)

本発明に係る文書類似性導出装置は、文からなる文書の文を形態素解析する手段と、形態素解析された文書から、当該文書に出現する索引語のTF/IDFによる重みを要素としたTF/IDFベクトルを求める手段と、形態素解析された文書から、当該文書に出現する体言について当該体言が出現する文中で共起した用言の頻度を要素とした共起ベクトルを求める手段と、形態素解析された文書から、当該文書に出現する索引語が出現する文章の文タイプを決定し、それぞれの文タイプの頻度を要素とした文タイプベクトルを求める手段とを含み、第1の文書のTF/IDF文書ベクトル、共起ベクトル及び文タイプベクトルを求め、第2の文書のTF/IDF文書ベクトル、共起ベクトル及び文タイプベクトルを求め、求めた第1の文書のTF/IDF文書ベクトル、共起ベクトル及び文タイプベクトルと第2の文書のTF/IDF文書ベクトル、共起ベクトル及び文タイプベクトルから第1の文書と第2の文書の類似性を求めるものである。

20

30

【0017】

このように本発明においては、TF-IDFベクトルに加え、少なくとも共起ベクトル及び文タイプベクトルの一方を用いて文書間の類似性を求めているので、より文書内の意味内容を反映した類似性を求めることができるという効果を奏する。

【0018】

文書は一以上の文からなる。したがって、第1の文書が1つの文からなり、第2の文書が複数の文からなる場合、第1の文書が複数の文からなり、第2の文書が1つの文からなる場合、第1の文書及び第2の文書ともに複数の文からなる場合、第1の文書及び第2の文書ともに1つの文からなる場合がある。つまり、本発明の文書類似性導出装置により、文間、複数の文からなる文書間、文-複数の文からなる文書間の類似性を求めることができる。

40

【0019】

共起ベクトルは、図2(上部)に示すように、文書に出現する体言について当該体言が出現する文中で共起した用言の頻度を算出し、その算出した頻度を要素として共起ベクトルを求める。図2(上部)では、例として、形態素解析により判明した「AAA」という体言について形態素解析により判明した「aa」という用言が3回文書中に出現したことを算出している。

【0020】

同様に、図2(下部)は文タイプベクトルも説示しており、文書に出現する索引語が出現する文章の文タイプを決定し、それぞれの文タイプの頻度を要素として文タイプベクトル

50

ルを求めている。図2(下部)では、例として、形態素解析により判明した「AAA」という体言について文タイプの決定処理を経て「QUESTION」の文タイプが3回文中に出現したことを算出している。

【0021】

(2) 文書群と文書の類似性(図3参照)

本発明に係る文書-文書群類似性導出装置は、前記文書類似性導出装置の各手段を含み、TF-IDFベクトルと共起ベクトル及び/又は文タイプベクトルである文書ベクトルを第1の文書について求め、複数文書からなる第2の文書群の各文書の文書ベクトルを求め、求めた第2の文書群の各文書の文書ベクトルから平均文書ベクトルを求め、求めた第2の文書群の平均文書ベクトルと第1の文書の文書ベクトルから第1の文書と第2の文書群の類似性を求めるものである。

10

このように本発明においては、文書間の類似性だけでなく、文書群と文書の類似性を求めることもできるという効果を有する。

【0022】

(3) 高い類似性を有した文書の特定

本発明に係る高類似性文書特定装置は、前記文書類似性導出装置の各手段を含み、TF-IDFベクトルと共起ベクトル及び/又は文タイプベクトルである文書ベクトルを比較対象となる比較対象文書について求め、第nの文書の索引TF-IDFの文書ベクトル、共起ベクトル及び文タイプベクトルを求め、比較対象文書の文書ベクトルと第nの文書の文書ベクトルから比較対象文書と第nの文書の類似性を求め、nは1ないしNまであり、各第nの文書と比較対象文書の類似性の中から類似性の高い第nの文書を特定するものである。

20

このように本発明においては、複数の文書と比較対象文書の類似性を求め、高い類似性を有する文書を特定するので、比較対象文書の内容によく類似した文書を得ることができるといふ効果を有する。

【0023】

(4) 高い類似性を有した文書群の特定

本発明に係る高類似性文書群特定装置は、前記文書類似性導出装置の各手段を含み、TF-IDFベクトルと共起ベクトル及び/又は文タイプベクトルである文書ベクトルを比較対象文書について求め、複数文書からなる第nの文書群の各文書の文書ベクトルを求め、求めた第nの文書群の各文書の文書ベクトルから平均文書ベクトルを求め、求めた第nの文書群の平均文書ベクトルと第1の文書の文書ベクトルから第1の文書と第nの文書群の類似性を求め、nは1ないしNまであり、各第nの文書群と比較対象文書の類似性の中から類似性の高い第nの文書群を特定するものである。

30

【0024】

(5) 回答支援システム(図4参照)

本発明に係る回答支援システムは、前記高類似性文書群特定装置を含み、前記各第nの文書群は類似する質問文からなり、比較対象文書も質問文であり、各第nの文書群の質問内容に対応する回答文を関連付けて予め記録し、前記高類似性文書群特定装置により類似性の高いとされた第nの文書群に関連付けられている回答文を出力するものである。

40

【0025】

後説する実施形態では、この回答支援システムを具体例として示したものである。特に、実施形態では、各第nの文書群を複数層のツリー構造にてデータベース化している。また、高類似性文書群特定装置は、このように各第nの文書群内の文書が相互に類似性が高くなるように、比較対象文書が属すべき最も類似性の高い第nの文書群を特定することにも用いることができ、それが後説する質問応答データベース構築支援システムとなる。

【0026】

図4は回答支援システムの発明原理図である。第1文書群から第N文書群までがあり、それぞれの文書群に対して予め平均文書ベクトルを求めて記録しておき、また、それぞれの文書群に対して共通の回答文を求めておく。そして、対象文書の文書ベクトルを求めて

50

、対象文書の文書ベクトルと記録している各文書群の平均文書ベクトルから対象文書と各文書との類似性を求め、最も高い類似性を有する第n文書群を特定し、この第n文書群の回答文を最適な回答文として使用者に出力する。なお、現在対象文書となっている文書も第n文書群に振り分けられ、新しく振り分けられた文書を含めて再度第n文書群の平均文書ベクトルを求める。新しく文書が振り分けられる度にしてもよいし、所定文書数蓄積された場合、所定期間毎に平均文書ベクトルを求めてもよい。同様に、新しい振り分けも所定文書数蓄積された場合、所定期間毎に実行してもよい。

【0027】

これまで装置又はシステムとして本発明を把握してきたが、所謂当業者であれば明らかであるように、プログラム又は方法としても把握することができる。

10

これら前記の発明の概要は、本発明に必須となる特徴を列挙したのではなく、これら複数の特徴のサブコンビネーションも発明となり得る。

【発明を実施するための最良の形態】

【0028】

[1. システム概要]

システムは、ユーザから質問メールを受け取ると、自動的に処理を始める。処理の結果、すなわち、回答候補は質問者への返答メールの形に整形されて最終査読者の査読を待つ。

【0029】

質問文の解析は、形態素解析のみを行い、その結果を言語データベースおよび知識データベースに照らして、質問応答データベース検索の準備をし、自動的に検索を実施する。言語データベースは、言語的素性を持つ辞書のようなものであり、知識データベースはパソコンに関する事典のようなものである。

20

質問文の解析として、形態素解析のみを活用する理由は、依存構造解析などの信頼性に起因するものである。

【0030】

最終査読者の査読済み質問応答データは、質問応答データベースの更新に供される。すなわち、質問応答データベースは査読済み回答文の返信をトリガとしてその質問応答データにより更新される。

【0031】

30

システム画面の例を図5に示す。いくつかの操作機能が用意されており、最終査読者が効率的な査読を行えるよう配慮されている。最終査読者がいくつかの回答候補のいずれかが妥当であると判断し所定のボタンを押下すると、その選択された回答候補が質問者に自動的にメール返信される。微細な加工が必要な場合には、この画面の上で加工し、それが質問者に回答メールとして自動的に返信される。つまり、最終査読者は選択した回答候補の文章を適宜修正し、メールを送信することができる。

【0032】

図5のシステム画面の構成は図6に説示する通りである。質問文入力フォームに質問文を入力し、カテゴリ判定ボタンを使用者が押下することで質問文と質問カテゴリの類似度が算出される。算出された類似度順にカテゴリのリストボックスに質問カテゴリを表示する。質問カテゴリ又は質問カテゴリに属する質問文を使用者が選択することで、テキストボックスに回答文が表示される。

40

【0033】

ここで、質問文入力フォームへの質問文の入力は、例えば、メーラの本文表示からテキストデータを貼り付けることで入力する。ただし、この例に限定されない。メーラに回答支援システムの機能を具備させてもよいし、逆に、回答支援システムにメーラの機能を具備させてもよい。さらには、メールシステムを用いることなく、質問者からの質問文を他の通信方法で取得する方式を適用することもできる。例えば、HTTP、FTPを用いることができる。

【0034】

50

[1.1 形態素解析]

文を形態素に分割して品詞を見分ける形態素解析については、自然言語処理の基礎技術の一つであり、所謂当業者であれば適宜適用が可能であるため、ここでは詳述しない。形態素解析エンジンとしては、例えば、MeCab、Chasen、KAKASIなどがある。

【 0035 】

[1.2 システム構成]

回答者が使用する回答者コンピュータ100、200上に回答支援システムを構築する。回答者コンピュータ100、200に回答支援プログラムがインストールされ、回答支援システムが構築される。本実施形態では、このように一つのコンピュータにより回答支援システムが構築されているが、クライアント・サーバ型で構築することもできる。例えば、クライアントでは、ユーザからの質問文をクライアントが受けてサーバに送信し、サーバで処理されて複数の回答候補をクライアントに返信する構成である。

10

【 0036 】

本実施形態の回答支援システムを構築したコンピュータの属するネットワーク構成の一例を図7に示す。LAN上に回答者コンピュータ100、回答者コンピュータ200、サーバ300、プリンタ(サーバ)400及びネットワーク機器500が接続され、相互に通信可能となっている。また、ネットワーク機器500は外部ネットワークとも接続し、他のコンピュータとLAN上のコンピュータを通信可能としている。ここでは、質問者であるユーザからメールが送信されるとして、メールサーバが送信するメールがネットワーク機器500を介して回答者コンピュータに送信される。回答者コンピュータが複数ある場合のメールの振り分け処理などは、コールセンタに構築されたシステムの一機能として実装され、周知・慣用技術であるためここでは詳述しない。

20

【 0037 】

回答支援システムが構築される回答者コンピュータ100は、例えば、CPU(Central Processing Unit)101、RAM102、ROM103、外部記憶装置であるHD(hard disk)104、CD-ROMからデータを読み出すCD-ROMドライブ105、入力装置であるマウス111及びキーボード112、出力装置であるディスプレイ121とスピーカー122、並びに、ネットワークに接続するためのLANインタフェース131からなる構成をとる。

30

回答者コンピュータ100の構成の一例を示したが、回答者コンピュータ200、サーバ300、ユーザコンピュータ600も同様の構成である。

【 0038 】

[2. 質問応答データベース]

質問応答データベースも回答者コンピュータ100、200にそれぞれ構築するものとする。ここで、別途データベースサーバとして構築し、複数の回答者コンピュータが共通に使用する構成にすることもできる。

【 0039 】

[2.1 質問応答データベースの構築手法]

質問応答メールデータとは、パソコンユーザから送られてきた問い合わせメールとそれに対する査読済み回答文のペアのことである。

40

質問応答メールデータの中には、同一データないしは類似データが多数存在する。したがって、問い合わせメールの内容または意味が同一または類似で、それらの回答文の文章表現も同じまたは類似である場合、それらを「類似データ」とみなす。

【 0040 】

メールコールセンターの質問応答データベースを構築するにあたって、実際には、10135件の質問応答メールデータを使用した。これらに対して、「質問カテゴリ」を作成し、類似データの「質問カテゴリ分類」を行った。質問カテゴリは、上に述べたように二層構造にし、上位層をブランチカテゴリ、下位層をリーフカテゴリと呼ぶ。類似データは同一リーフカテゴリに分類し、さらに相関関係があるリーフカテゴリは同一ブランチカテ

50

ゴリに分類する。

【 0 0 4 1 】

[2 . 2 質問応答データベースの構成]

構築したメールコールセンターの質問応答データベースは二段階（二層）のツリー構造である。リーフカテゴリはブランチカテゴリに属する場合もあるし、直接ルートカテゴリに属する場合もある。ルートカテゴリはブランチカテゴリの上位カテゴリであるが、ツリー構造のルートノードであるので、「層」とはみなさない。すなわち、全体を三層構造とは呼ばないことにしている。

【 0 0 4 2 】

ブランチカテゴリはデータを持たず、リーフカテゴリは同一データないしは類似データを持つ。ルートと各カテゴリの相関関係は下記の通りである：

- ・ルートカテゴリ ブランチカテゴリ リーフカテゴリ
- ・ルートカテゴリ リーフカテゴリ

この概略を図 8 に示す。

【 0 0 4 3 】

構築したメールコールセンターの質問応答データベースでは、10135 件の質問応答メールデータの内、利用対象外データ 3598 件を除き、計 6537 件に対して、83 個のブランチカテゴリおよび 634 個のリーフカテゴリが設定された。利用対象外データとは、いわゆるすなおな形・内容のものではなかったものなどであり、再利用にはむかないものなどである。説示中にでてきた数字はある検証実験で得られたものである。

【 0 0 4 4 】

[2 . 3 質問応答データベース構築支援システム]

質問応答データベースの構築には、多くの工数を要する。したがって、当初は、1705 の質問応答データについて人手で質問応答データベースを構築した。

質問応答データベース構築の効率をあげるため、その後、質問応答データベース構築支援システムを作成して活用した。上に述べた 6537 件の質問応答データは、この質問応答データベース構築支援システムを使用して構築したものである。そのスクリーンショットを図 9 に示す。なお、質問応答データベースは質問応答データベース構築支援システムを用いることなく、全て人手により構築してもよい。

【 0 0 4 5 】

[2 . 3 . 1 システム構成]

メールコールセンターの質問応答データベース構築支援システムの特徴は下記の通りである。

- ・分類する質問メールと既存の質問カテゴリの類似度を計算する
- ・操作しやすい G U I インタフェースを提供する
 - 質問メール、質問カテゴリに既存の質問メールの内容表示および質問カテゴリの表示
 - 分類する質問メールに対して、類似度順で質問カテゴリの提示
 - 分類する質問メールに対して、属する質問カテゴリの選択・作成・削除
 - 質問カテゴリごとの回答文テンプレートの作成(後記参照)

【 0 0 4 6 】

ここでは、メールコールセンターの質問応答データベース構築支援システムのシステム構成、システム用データベース、質問カテゴリ判定システムおよび G U I インタフェースについて述べる。

メールコールセンターの質問応答データベース構築支援システムの構成は図 10 で示す。

メールコールセンターの質問応答データベース構築支援システムでは、分類する質問メールが下記の 4 つのステップを通して質問カテゴリに分類される(回答文の作成に関しては後説)。

【 0 0 4 7 】

- 1) 分類する質問メールを質問カテゴリ判定システムを通して、既存の質問カテゴリと

10

20

30

40

50

の類似度を計算する。

【 0 0 4 8 】

2) GUI インタフェースで 1) で計算した類似度順ですべての既存の質問カテゴリを提示する。

【 0 0 4 9 】

3) 分類する質問メールが提示された既存のリーフカテゴリに属すると判断される場合、そのリーフカテゴリに分類する。自動的に処理することもできるが、本実施形態では GUI を介して使用者からの承認を経て分類している。具体的には 2) で提示したリスト形式で表示した質問カテゴリの使用者からの指定を受け付け、さらに、分類の承認を受け付ける。

10

【 0 0 5 0 】

4) 分類する質問メールが提示された既存のリーフカテゴリに属しないと判断される場合、分類する質問メールに対して、質問カテゴリの作成基準に従って、新しいリーフカテゴリまたはブランチカテゴリの作成を行う。分類する質問メールを新しく作成したリーフカテゴリに分類する。自動的に処理することもできるが、本実施形態では GUI を介して使用者からの承認を経て分類している。具体的には 2) で質問カテゴリが表示されない場合、質問カテゴリが表示された場合でも適当な質問カテゴリでないときに、使用者から新しいリーフカテゴリまたはブランチカテゴリの作成の指示を受け付ける。

【 0 0 5 1 】

[2 . 3 . 2 システム用データベース]

20

メールコールセンターの質問応答データベース構築支援システムで、質問メールの分類を行う際に、分類する質問メールデータを質問カテゴリ判定システムを通して、既存の質問カテゴリとの類似度の計算を行うため、事前に、メールコールセンターの質問応答データベース構築支援システム用の質問カテゴリを用意する必要がある。

ここで、メールコールセンターの質問応答データベース構築支援システム用データベースとして、1705 件の質問メールを利用して作成した質問カテゴリを利用する ([2 . 1 質問応答データベースの構築手法] を参照)。

【 0 0 5 2 】

[2 . 3 . 3 質問カテゴリ判定システム]

メールコールセンターの質問応答データベース構築支援システムでは、質問カテゴリ判定システムを利用して、質問メールが属するリーフカテゴリの判定を行う。

30

判定手法として、質問メールと質問カテゴリをベクトル空間上の点で表し、ベクトル間の類似度を定義する。

【 0 0 5 3 】

質問メールの文書ベクトルに関して、質問メールと質問メールが属する質問カテゴリの類似度が大きくなるように、ベクトルの要素を決定する。判定手法では、TF-IDF の重みづけによる文書ベクトルを拡張し、体言と用言の共起および文の特徴を考慮することで、質問メールの内容をより正確に反映する文書ベクトルを用いる。

【 0 0 5 4 】

質問メールの文書ベクトルは、下記の 3 種類である。

40

- ・ TF-IDF による文書ベクトル
- ・ 体言と用言の共起を考慮した文書ベクトル
- ・ 文タイプを考慮した文書ベクトル

また、質問カテゴリに属する質問メールの文書ベクトルを平均化したものを質問カテゴリの文書ベクトルとし、判定する質問メールの文書ベクトルとの重みづけ余弦尺度によって、両方の類似度を求める。類似度の計算結果によって、質問メールが属する質問カテゴリを判定する。

詳細には、後記 [3 . 特徴ベクトル] で説示する。つまり、特徴ベクトルは回答支援システムで質問文に対する適切な回答文を特定するだけでなく、質問応答データベース構築支援システムでも使用する。

50

【 0 0 5 5 】

[2 . 3 . 4 GUI]

メールコールセンターの質問応答データベース構築支援システムをツールとしてユーザが使用する際、容易に利用できるようにシステム用のGUIインタフェースを作成した。

図9で表示したボタンを押すことで、[2 . 3 . 1 システム構成] 冒頭で説明したインタフェース機能を実現することが可能である。

【 0 0 5 6 】

初期起動時、図9のウィンドウ左側に、分類するメール、リーフカテゴリ及びブランチカテゴリのリストを表示する。

分類する質問メール一件を選択してクリックすることで、図9のウィンドウ中央に選択した質問メールの内容が表示される。

10

【 0 0 5 7 】

図9で選択した質問メールに対して、「類似度計算」ボタンを押すことで、図9のウィンドウ左側のリーフカテゴリおよびブランチカテゴリが類似度順で再表示される。類似度順は色付で表示され、色が濃いほど類似度が高いことを示す。つまり、類似度を色の濃淡で顕示している。

【 0 0 5 8 】

リーフカテゴリをクリックして、リーフカテゴリに分類した質問メールの一覧が表示される。ここでリーフカテゴリを選択すると、そのリーフカテゴリに属する質問メールのリストが表示され、その中の質問メールをクリックすることで、図9のウィンドウ右側に質問メールの内容が表示され、使用者は参照することが可能である。ここで、分類する質問メールがどのリーフカテゴリに属するかを判断する。属するリーフカテゴリが存在する場合、使用者が「振り分け」ボタンを押すことで、属するリーフカテゴリに分類する。属するリーフカテゴリ存在しない場合、「(新)カテゴリ作成」ボタンを押して、属するリーフカテゴリを作成する。

20

【 0 0 5 9 】

[3 . カテゴリ回答文作成]

リーフカテゴリに分類された同一データないしは類似データの回答文は同じまたは類似であるため、リーフカテゴリに属する質問メールに対して、共通回答文であるリーフカテゴリ回答文を作成する。未知の問い合わせメールの回答文を作成する際、その質問メールの属するリーフカテゴリのカテゴリ回答文を用いて、回答文の作成を行う。

30

リーフカテゴリ回答文には、定型回答文とテンプレート回答文の2種類を用意する。定型回答文は機種関連情報などを含んでいない場合のためのものであり、回答文を作成する際にそのまま出力する。一方、テンプレート回答文は、ハードウェアやソフトウェアなどの多種類の機種関連情報に関するスロットが用意されており、それらの機種関連情報を機種関連情報データベースから抽出し、スロットに入れ、回答文を作成し出力する。

【 0 0 6 0 】

本章では、構築したメールコールセンターの質問応答データベースを利用し回答文の作成への応用について述べる。

メールコールセンターの質問応答データベースでは、リーフカテゴリは類似質問メールの集合である。類似質問メールとは、お問い合わせメールの内容または意味が類似して、そのお問い合わせメールに対して送信した最終査読データである回答文が類似または同じである質問メールのことを指す。

40

【 0 0 6 1 】

そのため、リーフカテゴリに属するすべての類似質問メールに対して、共通の回答文を持つと考えればよい。その共通の回答文はリーフカテゴリ回答文と定める。1個のリーフカテゴリに対して、1種類のリーフカテゴリ回答文を持つ。また、リーフカテゴリに属する類似質問メールを参照して、リーフカテゴリ回答文を作成することが可能である。

【 0 0 6 2 】

前説したようにリーフカテゴリ回答文は下記の2種類がある。リーフカテゴリ回答文に

50

OSや機種などの情報によって変更する部分を含んでいない定型回答文と、リーフカテゴリ回答文にOSや機種などの情報によって変更する部分を含んでいるテンプレート回答文である。

次節からそれぞれのリーフカテゴリ回答文について説明を行う。

【0063】

[3.1 回答文作成]

リーフカテゴリ回答文を利用して、お問い合わせメールの回答文を作成する方法について述べる。

お問い合わせメールの回答文を作成する流れを図11に示す。

お問い合わせメールは質問カテゴリ判定システムを通して、属するリーフカテゴリを判定する。属するリーフカテゴリを決定すれば、リーフカテゴリ回答文を利用して回答文を作成する。

10

【0064】

[3.1.1 定型回答文の利用]

お問い合わせメールが属するリーフカテゴリのリーフカテゴリ回答文が定型回答文である場合、そのリーフカテゴリ回答文をお問い合わせメールの回答文として出力する。

【0065】

[3.1.2 テンプレート回答文の利用]

お問い合わせメールが属するリーフカテゴリのリーフカテゴリ回答文がテンプレート回答文である場合、以下の4つのステップでお問い合わせメールの回答文を作成する。

20

【0066】

1) お問い合わせメールにある機種名を抽出する。

【0067】

2) 1)で抽出した機種名を利用して機種DBでテンプレート回答文にあるスロットに入る機種情報をマッチングする。

【0068】

3) 2)でマッチングした機種情報をスロットに入れ換える。

【0069】

4) 回答文を出力する。

【0070】

30

テンプレート回答文を利用してお問い合わせメールの回答文を作成する場合、事前に機種DBを作成しておく。本研究では、メールコールセンターでサポートする全機種に対して、テンプレート回答文を作成する際に利用したスロット項目に対応する機種情報を集めて、機種DBを作成した。

ここでは、機器名、機器情報について説示したが、当然他の情報であってもよい。

【0071】

[4.特徴ベクトル]

これより、分類済みの質問応答データベースを用いて、新たに入力として与えられた質問文がどのカテゴリに属するのかを判定する手法について説示する。

本実施形態では、質問文とカテゴリをベクトル空間上の点で表す。また、ベクトル間の類似度を定義する。質問文と、その質問文が属するカテゴリとの類似度が大きくなるようにベクトルの要素を決定し、類似度を定義することで、類似度によって質問文のカテゴリを推定するものである。

40

【0072】

本手法では、TF-IDFの重みづけによる文書ベクトルを拡張し、体言・用言の共起と、文の特徴を考慮することで、質問文の内容をより正確に反映した文書ベクトルを用いる。

つまり、以下の特徴ベクトルを複合的に用いる。

- ・TF (Term Frequency) / IDF (Inverse Document Frequency) による特徴ベクトル
- ・体言と用言の共起を考慮した特徴ベクトル

50

・文タイプを考慮した特徴ベクトル

【0073】

また、カテゴリに属する質問文の文書ベクトルを平均化したものをカテゴリの文書ベクトルとし、重み付きの余弦尺度によって類似度を求める。余弦尺度は、ふたつのベクトルの類似度を、ベクトルがなす角の余弦によって考えるもので、同じベクトル同士はそのなす角が0で余弦は1となり、完全に異なる要素を持つベクトル同士は直交して余弦は0になるというものである。ベクトル同士のなす角の余弦は以下の式で表せる。

【0074】

【数1】

$$\cos \vartheta = \frac{\vec{v} \cdot \vec{v}'}{|\vec{v}| \cdot |\vec{v}'|} \quad \cdot \cdot \cdot \quad (1)$$

($\vec{v} \cdot \vec{v}'$ は内積)

ベクトルが正規化済みであるならば、これはベクトル同士の内積に等しい。以降、基本的にベクトルはすべて正規化済みであると仮定する。つまり、余弦尺度は内積によって求められる。

【0075】

本実施形態においては、文書の特徴付ける語として、名詞(未知語を含む)、動詞、形容詞などの自立語の原型と品詞情報の組を用いることとする。語の解析は、日本語係り受け解析器Cabochaと、形態素解析器MeCabを用いて、形態素解析と、文節区切りまでを行った。半角と全角の同じ文字や、アルファベットの大文字小文字などを区別しないようにあらかじめ前処理した文を形態素解析器にかけ、連続する数字・アルファベット・記号はつなげて名詞とした。ただし、連続する名詞を複合名詞とすることは、学習データが不十分になると考え、行っていない。数詞と助数詞の連続については、数詞を実際の数字の並びではなく、数クラスに置き換える処理を行った。なお、括弧内の文で、二文節以上の文は、括弧内の文であるという情報は保持しつつ、別の一文として切り離して扱うようにしている。

【0076】

[4.1 TF/IDFによる特徴ベクトル]

システムでは、問い合わせメール中に出現する語のTF/IDFによる重みを要素とした文書ベクトルを拡張したもので質問文を表現する。

TF-IDF重み付けはテキストの自動索引づけにおいて、索引語の重みを計算する手法である。TF(Term Frequency)とは、ある文書dにおける索引語tの生起頻度であり、 $tf(d, t)$ と表記する。またIDF(Inverse Document Frequency)は文書の数Nと、索引語tが一回以上生起する文書の数 $dfreq(t)$ によって次のように定義される。

【0077】

【数2】

$$idf(t) = \log \frac{N}{dfreq(t)} + 1 \quad \cdot \cdot \cdot \quad (2)$$

索引語tの文書dにおける重み $w(t, d)$ として、TFとIDFの積をもちいるのがTF-IDF重み付けである。重み付けにTFを用いるのは、文書中で繰り返し生起する語はその文書において重要な概念であると考えられるためである。しかし、多くの文書に生起

10

20

30

40

50

する語は、文書を特定する性質を持たず、索引語として適していない。そこで、語がどのくらい特定性を持つかを I D F によって重み付けに反映させている。

【 0 0 7 8 】

[4 . 2 体言と用言の共起を考慮した特徴ベクトル]

T F / I D F による重み付けは、通常、ある語が特定の文書の特徴付ける尺度を表現するものであり、文の構造を反映しない。したがって、

- ・「電源を切る。」
- ・「電源を入れる。」

という二つの文に対して、「電源」という語は同じ重みが与えられる。だが実際には、目的とする質問文のカテゴリ判定においては、この二つは違う特徴を持つものとして認識すべきである。これは、語の出現頻度だけを考慮しては、とらえにくい特徴である。そこで、T F / I D F による重み付けに加えて、体言に対する用言の一文での共起の度合を重みとして用いることを考える。それぞれの体言について、一文中で共起した用言の頻度を要素とする特徴ベクトルを用いる。文書ベクトルの要素として、T F - I D F 重みと一緒に保持しておく。これにより、ふたつの語を比べた際に、共起ベクトルの余弦尺度による類似度を用いることを考える。

10

通常、T F - I D F のみによる文書ベクトル V と V' の類似度 $\text{sim}(V, V')$ は、余弦尺度、つまり内積によって求める。全文書中の語の数、すなわち文書ベクトルの次元を n とすると、以下のように表される。

【 0 0 7 9 】

20

【数 3】

$$\text{sim}(\vec{V}, \vec{V}') = \vec{V} \cdot \vec{V}' = \sum_{i=1}^n \{\omega_i \cdot \omega'_i\} \cdot \cdot \cdot \quad (3)$$

($\omega_i \cdot \omega'_i$ は TF-IDF による重み)

ここで、共起ベクトルの類似度を重みに加える。要素に TF-IDF による重みと体言・用言の共起ベクトルを持つ文書ベクトル V_c 、 V'_c の類似度 $\text{sim}(V_c, V'_c)$ を、以下のように定義する。

30

【 0 0 8 0 】

【数 4】

$$\text{sim}(\vec{V}_c, \vec{V}'_c) = \sum_{i=1}^n \{(\omega_i \cdot \omega'_i)(\vec{c}_i \cdot \vec{c}'_i)\} \cdot \cdot \cdot \quad (4)$$

($\vec{c}_i \cdot \vec{c}'_i$ は共起ベクトル)

40

上式は、ある語 i について、T F - I D F の重みが大きいほど、また、語 i に同じような共起の傾向があるほど、文書ベクトルの類似度が高くなる。上記の「電源」の例の場合、それぞれに共起している用言は「切る」、「入れる」であるので、共起ベクトルの類似度は 0 である。したがって文書ベクトル全体の類似度も 0 となり、ふたつの文は似ていないと判断される。

【 0 0 8 1 】

[4 . 3 文タイプを考慮した特徴ベクトル]

パソコンユーザから送られてきた問い合わせメールの内容をより正確に反映した特徴ベクトルを作成するため、文中の語がどのような意味の文に出現するのか、という傾向につ

50

いて考える。そのために、まず問い合わせメールを分析してそれぞれの文タイプごとの特徴を調べ、分析結果をもとに文タイプ同定のルールを作成する。

【 0 0 8 2 】

[4 . 3 . 1 質問文の分析]

ここでは、質問メールを分析することにより、次のように少数の文タイプを設定した。

- ・ Question : 「～できますか?」「～を教えてください」など、質問を述べてある文。
- ・ Problem : 「～ができません」「～する方法がわかりません」など、問題を述べてある文。
- ・ Intention : 「～したい」「～しようと思う」など、質問者の意図・希望が述べてある文。
- ・ Situation : 問題発生の手順・状況などについて述べてある文。
- ・ Think : 「～だと思えます」など、質問者の考えが述べてある文。
- ・ Other case : 「HDDでの再生は問題ありません」など、別の状況では問題が発生しない場合が述べてある文。
- ・ About : 「～について」などの、質問内容を端的に表している文。質問、回答の一行目に述べられることがある。
- ・ Message : エラーメッセージや、ダイアログなど、画面に表示された文字列の内容を述べてある文。
- ・ etc : その他の情報

10

【 0 0 8 3 】

[4 . 3 . 2 質問文の分析結果]

上記の文タイプを集計した結果を以下に示す。

- ・ Question 324
- ・ Problem 648
- ・ Intention 87
- ・ Situation 398
- ・ Think 37
- ・ Other case 80
- ・ About 368
- ・ Message 96
- ・ etc 34

20

30

【 0 0 8 4 】

質問について述べてある文や、パソコンの不具合・問題について述べてある文など、上記9種類の文タイプを設定し、約一週間分の問い合わせメール、323件2072文を分析して、文末表現や機能語から、文タイプを同定するルールを作成し、各文タイプ中での各単語の頻度を要素とする特徴ベクトルを用いる。上の表の右端の数字は、参考までに記したものであり、それぞれの文タイプについて、分析の際に出現した回数である。

【 0 0 8 5 】

ほとんどの質問には、QUESTIONかPROBLEMのどちらかが含まれ、どちらも出現しない場合は323件中に3件だけであった。その3件中のすべてにINTENTIONが含まれていた。さらに、それぞれのタイプについての分析を以下に示す。

40

- ・ Question : ほとんどの場合文末が記号「?」か助詞「か」、あるいは「教えてください」「ご教示ください」「お願いします」などで終わる。その他のタイプはほとんどマッチしない。
- ・ Problem : 文末が自立の動詞・形容詞の基本形や、「～できません」「～しない」「～してしまう」などで終わる場合が全体の3分の2を占める。また、QUESTION文の直前に多く出現する。
- ・ Intention : ほとんどの場合、文末が「～したい」「～ほしい」「～しようと思っています」などで終わる。
- ・ Situation : 「～しました」のような過去形で終わる場合が多いが、そうでない場合も

50

多くある。PROBLEM文の直前に多く出現する。

- ・Think : 「～かと」を含むか、文末が「思う」「気がする」などで終わる。
- ・Other case : 「～は」「～では」「～も」「～と」「～だと」などを含む文で、文末が「できる」「異常ない」「問題ない」「正常です」「発生しない」などで終わる。
- ・About : 質問、回答の一行目において、文末が名詞で終わる。
- ・Message : 文の全部、あるいは一部が「」や'で括られていることが多い。その直後に「という」「と、」「って」などの語がつき、「表示されました」「出ました」「メッセージが出ました」などの文が続く。
- ・etc : 「初心者です」「名前は～です」などの情報がある。これらについては、あらかじめ対応ルールを用意しておくのが難しく、また出現頻度も少ないため、今回は対応を見送ることにする。

10

【0086】

[4.3.3 文タイプ同定ルール]

上記の分析結果をもとに、文タイプの同定ルールを作成した。ルールは三段階に分けて適用される。まず、最初に適用するルールについて述べる。以下のそれぞれにあてはまる文に、重複を許してタイプを割り振っていく。

上記の分析結果をもとに、文タイプの同定ルールを作成した。ルールは三段階に分けて適用される。まず、最初に適用するルールについて述べる。以下のそれぞれにあてはまる文に、重複を許してタイプを割り振っていく。

【0087】

20

- ・Question : 文末が「が?」以外の疑問符で終わる。あるいは、文末が助詞「か」で終わる。あるいは、文末の5文節以内に「教えて」「教示」「教授」「お知らせ」「なぜ」「願い」を含む。
- ・Problem : 文末が自立の動詞・形容詞の基本形で終わる。あるいは、文末が「でした」「が」でなく、格助詞「が」を含む文のうち、格助詞「が」と文末の間に他の助詞を含まない。あるいは、文末の3文節に「すみません」「すいません」「していません」「してありません」を含まず、「なくなっています」「なくなった」「なくなり」「てしまった」「なります」「なりました」「まいりました」「まいります」「ません」「ない」「しまう」「れる」「れます」を含む。
- ・Intention : 文末の5文節に「(動詞)+たい」「ほしい」「(動詞)+(よ)うと」を含み、その後動詞の「思う」「考える」が続く。
- ・Think : 文中に助詞の並び「かと」を含む、あるいは、文末の3文節に「思う」「思った」「思われ」「考えられ」「気がする」「気がします」を含む。
- ・Other case : 助詞、あるいは助詞の並び「は」「では」「も」「と」「だと」を含む文で、文末が「できる」「動く」「作動(する)」「動作(する)」「起動(する)」の活用のうち、「基本形」「た」「ます」「ている」で終わるか、「異常」「問題」「不都合」の後に「ありません」「なかった」「ない」が続いて終わる。あるいは、助詞「は」「と」の後に、「正常に」「正しく」「普通に」「通常」「きちんと」「うまく」「ちゃんと」を含む文がくる。

30

・About : 質問、回答の一行目において、文末が名詞で終わる。

40

・Message : 助詞「と」を含む文で、以降に「メッセージ」「ボックス」「ポップアップ」「表示」「エラー」動詞「出る」が出現する。助詞「と」の直前に、「」、()、'、`、`、で括られた部分がある場合、複数の文にまたがっている場合でも、括弧などで括られた内部を全て`MESSAGE`と判断する。

【0088】

次に、なにもタイプが割り振られなかった文に対して、SITUATIONか、PROBLEMを割り振る。次の三種類の場合を考える。

・質問文中にQUESTIONもPROBLEMも出現していない場合。ABOUTが出現している場合は、ABOUT中の語を含む文をPROBLEMとする。ABOUT中の語を含む文がない場合や、ABOUTが出現していない場合、タイプが割り振られていない一番最初の文をPROBLEMとする。残りはSITUA

50

TIONとする。

・質問文中にQUESTIONが出現している場合。QUESTIONの直前の文にタイプが割り振られていない場合、PROBLEMとする。残りはSITUATIONとする。

・それ以外の場合。タイプが割り振られていない文をすべてSITUATIONとする。

【0089】

次に、重複した文タイプに対して、タイプ間の優先順位にもとづいたルールを適用してタイプを確定する。

文タイプが重複している文は、以下の優先順位で文タイプを決定する。

``MESSAGE'' > ``ABOUT'' > ``QUESTION'' > ``PROBLEM'' > ``OTHERCASE'' > ``INTENTION'' > ``THINK''

10

【0090】

以上のルールを適用して、文タイプを決定する。質問文中に出現する語は、どのような文タイプ中で何度出現するのかという情報をベクトルとして持つことになる。

求めた文タイプを利用して、式(4)を次のように拡張する。要素にTF-IDFによる重み、体言・用言の共起ベクトル、文タイプベクトルを持つ文書ベクトル V_t, V'_t の類似度 $sim(V_t, V'_t)$ を、以下のように定義する。

【0091】

【数5】

$$sim(\vec{V}_t, \vec{V}'_t) = \sum_{i=1}^n \{(\omega_i \cdot \omega'_i)(\vec{c}_i \cdot \vec{c}'_i)(\vec{t}_i \cdot \vec{t}'_i)\} \quad (\vec{t}_i, \vec{t}'_i \text{ は文タイプベクトル})$$

20

$$sim(\vec{V}_t, \vec{V}'_t) = \sum_{i=1}^n \{(\omega_i \cdot \omega'_i)(\vec{c}_i \cdot \vec{c}'_i)(\vec{t}_i \cdot \vec{t}'_i)\} \cdot \cdot \cdot \quad (5)$$

(\vec{t}_i, \vec{t}'_i は文タイプベクトル)

式(5)は、語iが同じような文タイプに出現する傾向がある場合、文書ベクトルの類似度が大きくなることを表している。質問文が、どのようなことについて述べているのか、という傾向が似ているものを、類似度が高い、と評価する。

30

【0092】

[4.4 カテゴリの平均ベクトルを用いた類似度計算]

未知の問い合わせメールがどの質問カテゴリに属するのかを計算するのに、各質問カテゴリ内の質問文の特徴ベクトルを平均化したものを便宜的に質問カテゴリのベクトルとする。その概念図を図12に示す。

これらに対して、未知の問い合わせメールのベクトルとの類似度を計算する。これは、カテゴリ内の質問文の文書ベクトルを平均化することで、少数のノイズを取り除き、カテゴリ内で真に特徴的な語の情報のみを残すことができるからである。

【0093】

【数6】

40

$$\vec{A} = \frac{\sum_{i=1}^n \vec{a}_i}{n} \cdot \cdot \cdot \quad (6)$$

\vec{A} : n個の正規化前の質問文の文書ベクトル $\vec{a}_1 \dots \vec{a}_n$ を持つカテゴリ

Aの平均文書ベクトル

50

そして、質問ベクトルと各質問カテゴリとの距離を計算し、最も近い質問カテゴリをもつカテゴリに質問が属すると判断する。

この平均ベクトルを正規化したものと、未知の質問文の文書ベクトルとの類似度の計算結果を利用して、回答作成支援システムを作成する。

【0094】

[4.5 具体例]

[4.5.1 文書ベクトル]

W: 単語空間

W_i: ある体言と対応している

例: $i = PC$, $i' = 電源$. . .

ここにおいて、ある文書ベクトルV内の単語iをTF-IDFにおいて重み付けした値をTF-IDF(i)と表すものとする。

【0095】

この時、例として下に挙げる文章1、2における文章ベクトルは文中の体言iを軸とし

$w: TF-IDF(i)$

$c: \{v: TF-IDF(v), v': TF-IDF(v'), \dots\}$ v, v': 体言iと文中で共起する用言、 $c: v$ を軸として持つベクトル

$t: \{文タイプ1: 文タイプ1中でのiの出現回数, \dots\}$ t: 文タイプを軸として持つベクトル

以上の3つの値をセットそして持ち、そのベクトルの要素wを正規化したものとする。

【0096】

文章1

買って来たばかりのPCの電源が入りません。

何をしたらいいのでしょうか？

[文章ベクトルV]

$\{i_1: [w = TF-IDF(PC), c = \{買う: TF-IDF(買う), 入る: TF-IDF(入る)\}, t = \{否定: 1\}]\}$

$i_2: [w = TF-IDF(電源), c = \{買う: TF-IDF(買う), 入る: TF-IDF(入る)\}, t = \{否定: 1\}]\}$

$i_3: [w = TF-IDF(何), c = \{する: TF-IDF(する)\}, t = \{疑問: 1\}]\}$ * $i_1 = PC$, $i_2 = 電源$, $i_3 = 何$

【0097】

文章2

PCが起動しないのですが、どうしたらよろしいですか？

[文章ベクトルV']

$\{i_1: [w = TF-IDF(PC), c = \{TF-IDF(起動)\}], t = \{疑問: 1\}]\}$ * $i_1 = PC$

【0098】

[4.5.2 類似度]

文章ベクトルV、V'の類似度を計算しようとする時、ベクトル空間の次元数は(V、V')の単語空間の次元数に等しい。

よって、例におけるVは軸としてPCしか持たない1次元のベクトル空間であるがここでは(V、V')の単語空間に拡張する。

また共起ベクトルc、文タイプベクトルtも同様に拡張した文章ベクトルV'をV''と表すとそれは以下のようなになる。

【0099】

[文章ベクトルV'']

$\{i_1: [w = TF-IDF(PC), c = \{買う: TF-IDF(買う), 入る: TF-IDF(入る), 起動: TF-IDF(起動), する: TF-IDF(する)\}, t = \{疑問: 1, 否定: 0\}]\}$

10

20

30

40

50

i_2 : [$w = TF - IDF$ (電源), $c = \{ \text{買う} : TF - IDF$ (買う), $\text{入る} : TF - IDF$ (入る), $\text{起動} : TF - IDF$ (起動), $\text{する} : TF - IDF$ (する) } , $t = \{ \text{疑問} : 0$, $\text{否定} : 0 \}$]

i_3 : [$w = TF - IDF$ (何), $c = \{ \text{買う} : TF - IDF$ (買う), $\text{入る} : TF - IDF$ (入る), $\text{起動} : TF - IDF$ (起動), $\text{する} : TF - IDF$ (する) } , $t = \{ \text{疑問} : 0$, $\text{否定} : 0 \}$] } * $i_1 = PC$, $i_2 = \text{電源}$, $i_3 = \text{何}$

ここにおいて類似度の計算は以下の式(5)に従う。また、ここにおける($t \cdot t'$)は文タイプが一致すれば1一致しないならば0を返すものである。

【0100】

[5.動作]

図13は本実施形態に係る回答支援システムのブロック構成であり、図14は本実施形態に係る回答支援システムの動作フローチャートである。なお、図13に示したブロック構成は一例であり、所謂当業者で明らかであるように複数のモジュール構成をとることができる。そして、ここでは、動作主体を明示しているが、ハードウェア的視点から言えば、コンピュータ、プロセッサが動作主体である。

【0101】

質問文が内包された質問メールを質問者がユーザコンピュータ600で作成し、回答者コンピュータ100のアドレス宛に送信する。

回答者コンピュータ100は複数のメールサーバを介してユーザコンピュータ600からの質問メールをメーラで受信する。なお、回答者コンピュータ100が直接アクセスするメールサーバが所定メールアドレスのメールを、登録された回答者コンピュータへ適宜振り分けする機能を有する構成であってもよい。

【0102】

使用者はメーラで受信した質問メールを本回答支援システムへ取り込む指示を行う。

入力部1は使用者から指示を受け付け、指示された質問メールを取り込む(S100)。

前処理部2は全角(半角)文字変換やアルファベットの大文字(小文字)変換などの前処理を実行する。

形態素解析部3は前処理後の質問文を形態素解析する(S200)。

【0103】

文書ベクトル作成部4のTF-IDF文書ベクトル部41、共起ベクトル部42及び文タイプ文書ベクトル部43はそれぞれ取り込んだ質問メールの本文の各ベクトルを求める。

類似性算出部6は、各質問カテゴリの平均文書ベクトルを読み出し、この読み出した平均文書ベクトルと求めた質問メールの各ベクトルから式(5)を用いて類似度を求める(S300、S400)。

類似性算出部6が各質問カテゴリとの類似度を求めた後に、出力部8は各質問カテゴリを読み出し、類似度順に質問カテゴリをリスト表示する(S500)。

【0104】

使用者からの質問カテゴリの選択を受け付け、回答文書特定部7が質問カテゴリの識別情報から質問カテゴリと関連付いて記録している回答文を読み出し、出力部8がそれを表示する(S600)。

【0105】

出力部8が表示している回答文への修正を受け付ける。

使用者からの承認を受け付けると、メーラを介して回答文が質問者に返信される(S700)。

【0106】

さらに、質問応答データベース構築支援システムを介して今回送信されてきた質問文、その回答文及び平均文書ベクトルが記録される(S800)。ここで、既に質問カテゴリも決定されており、使用者から質問カテゴリの選択を受け付けることなく迅速に記録処理

10

20

30

40

50

がなされる。

【実施例】

【0107】

カテゴリ判定の精度を確かめるために、3種類の実験を行い、結果の評価する。

(カテゴリ判定実験)

実験データとして、分類済みの質問・回答データのうち、1カテゴリに3件以上の質問文を持つ629カテゴリをデータAとして用いる。また、1カテゴリに12件以上の質問文を持つ145カテゴリをデータBとして用いる。データAの総データ数は6536件で、83個の上位カテゴリを持つ。データBの総データ数は4023件で、52個の上位カテゴリを持つ。これらのデータに対し、データを3分割してそのうちふたつを学習データとして使い、残りをテストデータとして3回テストを行った結果の平均をとる3分割交差検定を行い、質問文の正解カテゴリと、正解カテゴリの上位カテゴリを何位に判定したかを調べた。

10

質問文に対して、正解カテゴリと、正解の上位カテゴリをどれだけ上位に判定したかを評価とする。

【0108】

【表1】

| データ | 一位に判定 | 三位以内 | 上位カテゴリを一位に判定 | 上位カテゴリを三位以内 |
|----------|-------------|-------------|--------------|-------------|
| A(6536件) | 3043(46.6%) | 4562(69.8%) | 4463(68.3%) | 5798(88.7%) |
| B(4023件) | 2505(62.3%) | 3459(86.0%) | 3202(79.6%) | 3805(94.6%) |

20

カテゴリ判定の実験結果である。

データAでの判定結果は、データBに比べると悪い。その理由として、学習データの不足が考えられる。データAのカテゴリ数は629個であるが、その大半が1カテゴリ内に3個か4個程度のデータしか持っていない。学習データの数に比べて、カテゴリ数が非常に多いため、判定ミスが増加したものと考えられる。

30

【0109】

(TF-IDFによる重み付けと、提案手法との比較実験)

判定実験での実験データBを用いて、単純なTF-IDFによる重み付けだけを用いる文書ベクトルと、提案手法である、体言・用言の共起と、文タイプを考慮した文書ベクトルの類似度による判定精度を比較した。

質問文に対して、正解カテゴリと、正解の上位カテゴリをどれだけ上位に判定したかを評価とする。

【0110】

【表2】

実験データB(総データ数4023件)による結果

| 文書ベクトル | 一位に判定 | 三位以内 | 上位カテゴリを一位に判定 | 上位カテゴリを三位以内 |
|--------|-------------|-------------|--------------|-------------|
| TF-IDF | 2343(58.2%) | 3342(83.1%) | 3148(78.3%) | 3768(93.7%) |
| 提案手法 | 2505(62.3%) | 3459(86.0%) | 3202(79.6%) | 3805(94.6%) |

40

【0111】

TF-IDFと提案手法の文書ベクトルによる比較実験の結果である。

50

提案手法の方が、若干精度が良いが、改善率としては一割程度である。その理由として、ベクトルを平均化した結果、カテゴリの文書ベクトルが非常に特徴的になってしまった。カテゴリの文書ベクトル同士の類似度を計算したところ、ほとんど0に近い値ばかりになった。特徴的な語があると、それに強く反応してしまい、他の要素があまり考慮されていない。

【0112】

(KNN法と平均ベクトル法によるカテゴリ判定の比較実験)

判定実験での実験データBを用いて、KNN法によるカテゴリ判定と、平均ベクトル法によるカテゴリ判定の精度を比較した。平均ベクトル法では、未知の質問文の文書ベクトルを入力として、カテゴリの平均ベクトルとの類似度を用いてカテゴリを判定したが、KNN法では入力ベクトルとすべての学習データ内の文書ベクトルとの類似度を求め、類似度が高い方からk個の文書ベクトルが属するカテゴリから、入力ベクトルの属するカテゴリを判定する。

10

【0113】

データBでは、学習データ内のすべてのカテゴリが最低8個のデータを持っていることが保証されているので、kの値は8とした。

類似度の計算は、提案手法である体言・用言の共起と文タイプを考慮した文書ベクトルを用いて計算した。

質問文に対して、正解カテゴリと、正解の上位カテゴリをどれだけ上位に判定したかを評価とする。

20

【0114】

【表3】

実験データB(総データ数4023件)による結果

| | 一位に判定 | 三位以内 | 上位カテゴリを一位に判定 | 上位カテゴリを三位以内 |
|---------|-------------|-------------|--------------|-------------|
| KNN法 | 2402(59.7%) | 3194(79.4%) | 3243(80.6%) | 3768(92.7%) |
| 平均ベクトル法 | 2505(62.3%) | 3459(86.0%) | 3202(79.6%) | 3805(94.6%) |

30

【0115】

KNN法と平均ベクトル法の比較実験の結果である。

上位カテゴリの一位判定において、KNN法のほうがわずかに高い値を出しているものの、全体的には、特に三位以内での判定において、平均ベクトル法のほうが良い精度を出している。

KNN法の判定ミスの理由として、類似度の高いk個のデータの中に、正解のカテゴリに属するデータがひとつもない場合が4023件中583件もあることが挙げられ、データの分布がかなりの範囲で重なっている。

【0116】

以上の前記各実施形態により本発明を説明したが、本発明の技術的範囲は実施形態に記載の範囲には限定されず、これら各実施形態に多様な変更又は改良を加えることが可能である。そして、かような変更又は改良を加えた実施の形態も本発明の技術的範囲に含まれる。このことは、特許請求の範囲及び課題を解決する手段からも明らかなことである。

40

【図面の簡単な説明】

【0117】

【図1】発明の原理図(文書の類似性)である。

【図2】発明の原理図(共起ベクトル、文タイプベクトル)である。

【図3】発明の原理図(文書-文書群の類似性)である。

【図4】発明の原理図(回答支援)である。

【図5】本発明の実施形態に係るシステム画面である。

50

【図6】図5のシステム画面の構成である。

【図7】本発明の実施形態に係る回答支援システムを構築したコンピュータの属するネットワーク構成の一例である。

【図8】本発明の実施形態に係る質問応答データベースのツリー構造である。

【図9】本発明の実施形態に係る質問応答データベース構築支援システムのスクリーンショットである。

【図10】本発明の実施形態に係るメールコールセンターの質問応答データベース構築支援システムの構成である。

【図11】本発明の実施形態に係る質問メールの回答文を作成する動作フローチャートである。

10

【図12】本発明の実施形態に係る各質問カテゴリ内のベクトル概念図である。

【図13】本発明の実施形態に係る回答支援システムのブロック構成である。

【図14】本発明の実施形態に係る回答支援システムの動作フローチャートである。

【符号の説明】

【0118】

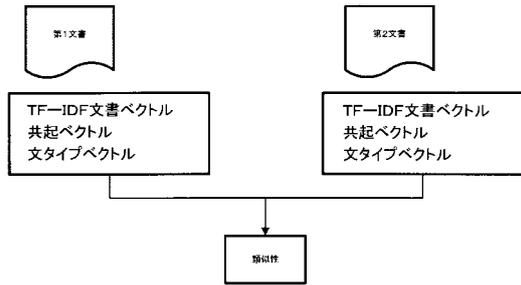
- 1 入力部
- 2 前処理部
- 3 形態素解析部
- 4 文書ベクトル作成部
 - 41 TF-IDF文書ベクトル部
 - 42 共起ベクトル部
 - 43 文タイプ文書ベクトル部
- 5 文書ベクトル記憶部
- 6 類似性算出部
- 7 回答文書特定部
- 8 出力部
 - 100 回答者コンピュータ
 - 101 CPU
 - 102 RAM
 - 103 ROM
 - 104 HD
 - 105 CD-ROMドライブ
 - 111 マウス
 - 112 キーボード
 - 121 ディスプレイ
 - 122 スピーカー
 - 131 LANインタフェース
- 200 回答者コンピュータ
- 300 サーバ
- 400 プリンタ
- 500 ネットワーク機器
- 600 ユーザコンピュータ

20

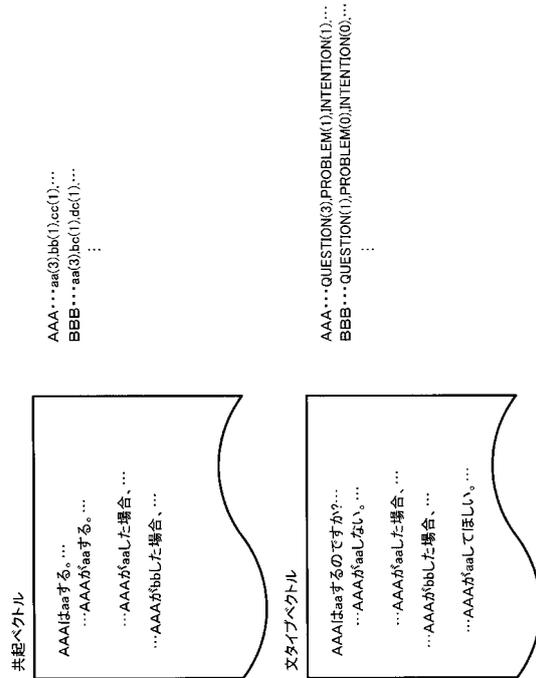
30

40

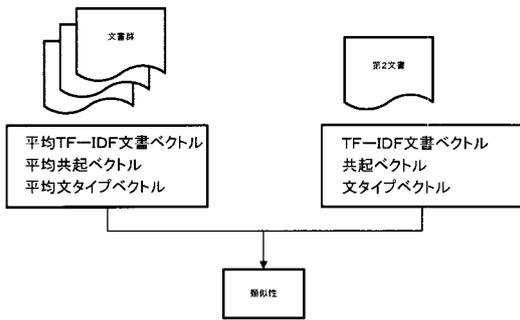
【図1】



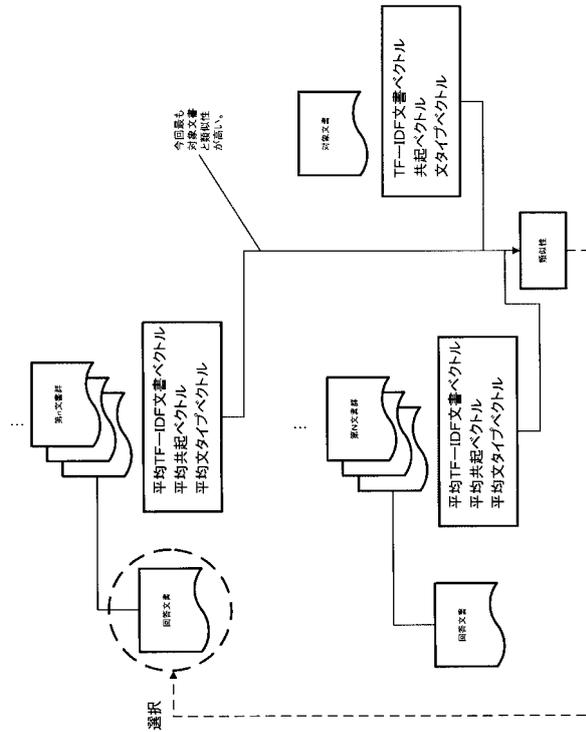
【図2】



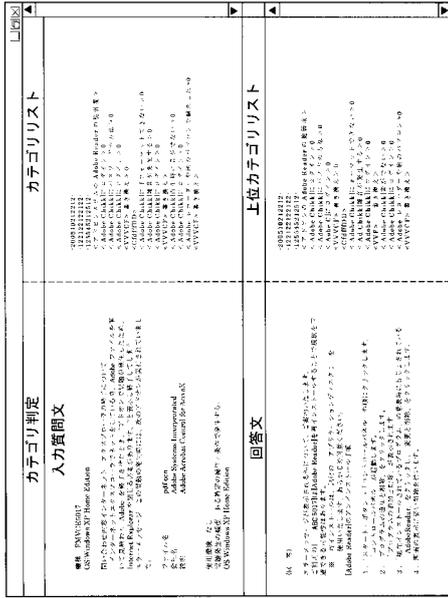
【図3】



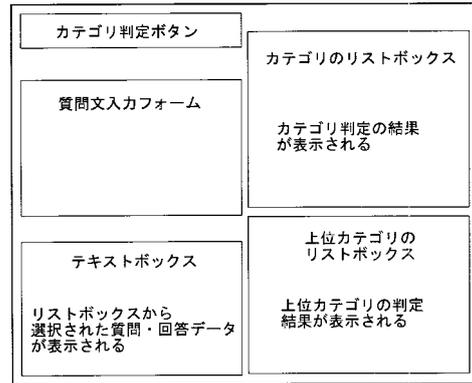
【図4】



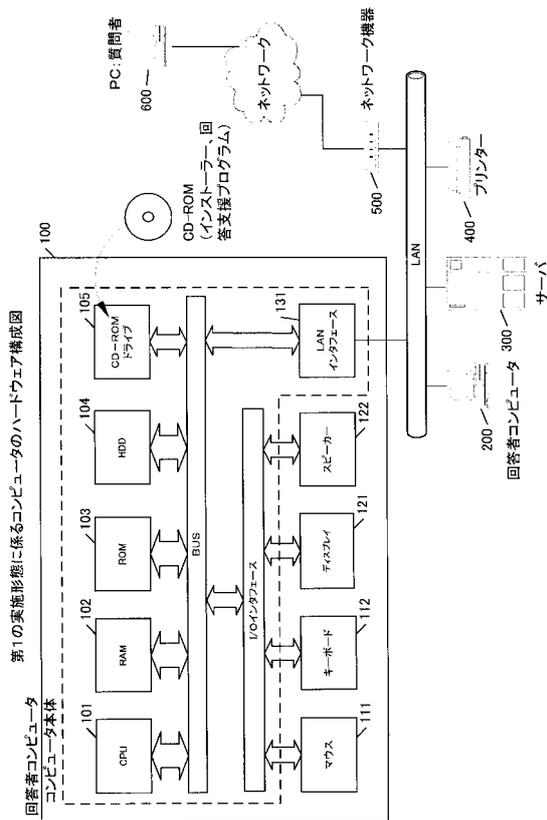
【図5】



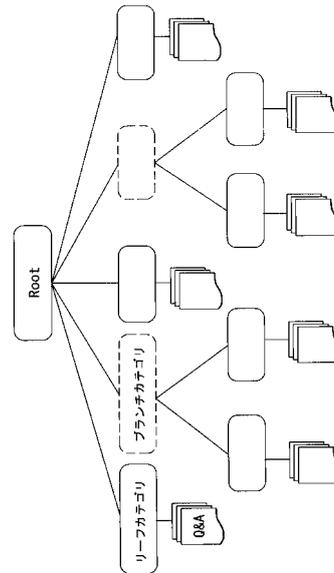
【図6】



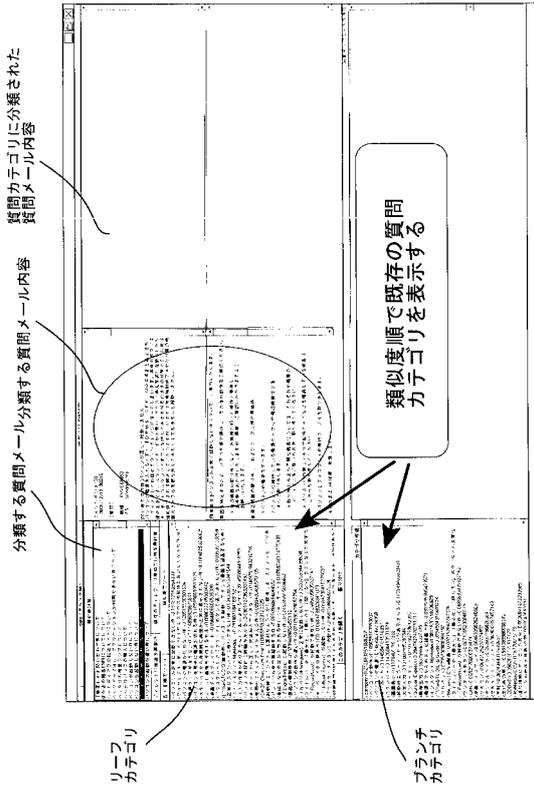
【図7】



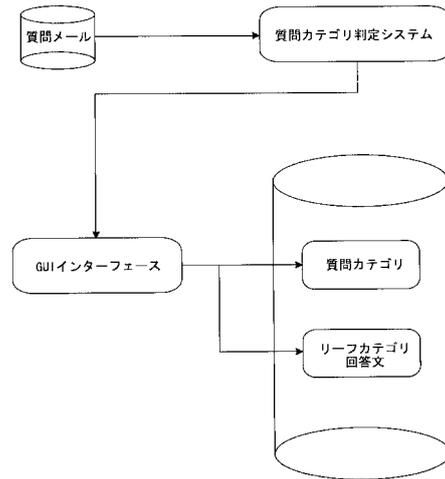
【図8】



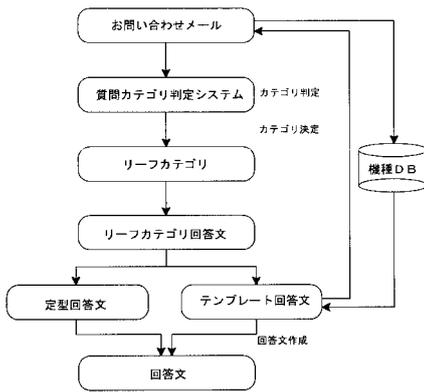
【図9】



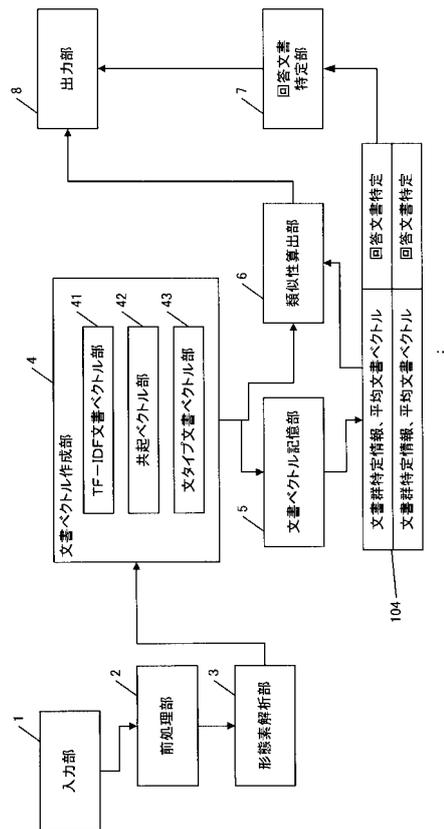
【図10】



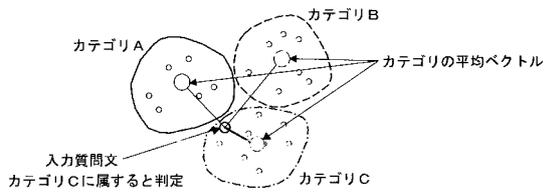
【図11】



【図13】

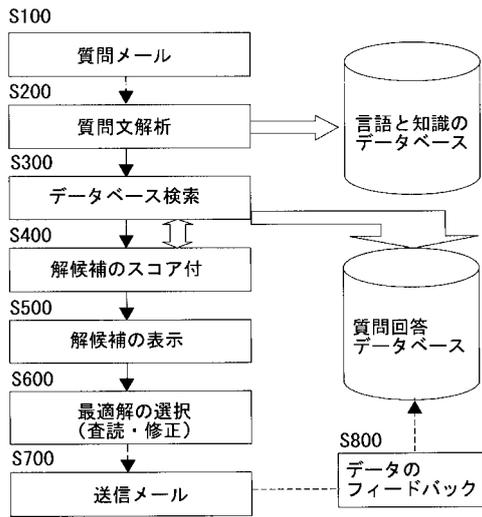


【図12】



【図14】

システム構成と処理の概要



フロントページの続き

- (56)参考文献 特開2002-245067(JP,A)
特開2001-331515(JP,A)
特開平10-078971(JP,A)
特開平06-110929(JP,A)
大森信行、外3名、tf・idf法を用いた関連マニュアル群のハイパーテキスト化、情報処理学会研究報告(97-NL-121)、日本、社団法人情報処理学会、1997年 9月12日、第97巻、第85号、p.111-118
岩崎礼次郎、外1名、コールセンターにおける対話データを用いた営業日報の自動生成、第22回 ことば工学研究会資料(SIG-LSA503)、日本、社団法人人工知能学会 ことば工学事務局、2006年 3月10日、p.87-94
岡村潤、外3名、複数マニュアルの自動ハイパーテキスト化における類似度計算手法について、情報処理学会研究報告(98-FI-51)、日本、社団法人情報処理学会、1998年 9月18日、第98巻、第81号、p.71-78
高木徹、外1名、単語出現共起関係を用いた文書重要度付与の検討、情報処理学会研究報告(96-FI-41)、日本、社団法人情報処理学会、1996年 4月18日、第96巻、第34号、p.61-68

(58)調査した分野(Int.Cl., DB名)

G06F 17/30