

(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2004-118647

(P2004-118647A)

(43) 公開日 平成16年4月15日(2004.4.15)

(51) Int. Cl.<sup>7</sup>  
G06F 17/30

F I

G06F 17/30 210B  
G06F 17/30 170A  
G06F 17/30 180A  
G06F 17/30 330C  
G06F 17/30 340B

テーマコード(参考)

5B075

審査請求 有 請求項の数 16 O L (全 27 頁)

(21) 出願番号 特願2002-282795 (P2002-282795)  
(22) 出願日 平成14年9月27日(2002.9.27)

(71) 出願人 301022471  
独立行政法人通信総合研究所  
東京都小金井市貫井北町4-2-1  
(74) 代理人 100121511  
弁理士 小田 直  
(74) 代理人 100097836  
弁理士 福井 國敏  
(72) 発明者 村田 真樹  
東京都小金井市貫井北町4-2-1 独立  
行政法人通信総合研究所内  
Fターム(参考) 5B075 ND03 NK32 NK33 PP24 PR04  
QP01 QP03

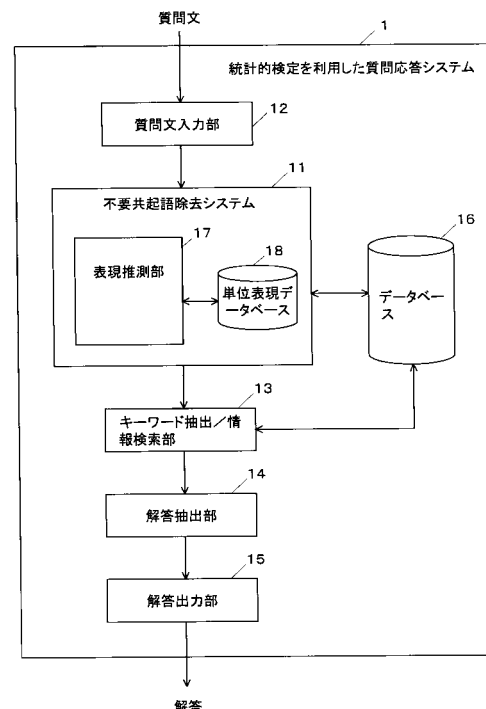
(54) 【発明の名称】 統計的検定を利用した質問応答方法、質問応答システム、質問応答プログラムおよび質問応答プログラムを記録した記録媒体

(57) 【要約】

【課題】 質問文に対する解が数量表現の場合に、当該数量表現に付される単位表現として適切な単位表現かを的確に判断することが可能な質問応答方法および質問応答システムを提供する。

【解決手段】 入力された質問文情報とデータベース16に記憶された電子化テキスト情報とに基づいて、不要共起語削除システム11の表現推測部17において、統計的検定を用いて、前記電子化テキスト情報において出現する前記質問文の主たる名詞と数量表現との組み合わせパターンから抽出される単位表現のうち前記質問文中の主たる名詞と共起して出現する可能性が低い単位表現を不要な単位表現と判断して除去する。

【選択図】 図1



## 【特許請求の範囲】

## 【請求項1】

自然言語による質問文を入力し、データベースに記憶された電子化テキスト情報中の文との照合によって解を生成して出力する質問応答方法であって、  
前記質問文を入力する質問文入力過程と、  
入力された質問文の解が数量表現である場合に、統計的検定を用いて、前記電子化テキスト情報において出現する前記質問文の主たる名詞と数量表現との組み合わせパターンから抽出される単位表現のうち前記質問文中の主たる名詞と共起して出現する可能性が低い単位表現を不要な共起語と判断し、前記質問文中の主たる名詞と共起して出現する可能性が高い単位表現を有用な共起語と判断し、前記有用な共起語と判断した単位表現を用いて前記質問文の解としての数量表現を推測する表現推測過程と、  
前記推測された数量表現を用いて前記質問文の解を抽出する解答抽出過程とを有することを特徴とする統計的検定を利用した質問応答方法。

10

## 【請求項2】

前記表現推測過程は、  
前記質問文情報に基づいて、質問文の解が数量表現であるかを判断する数量表現判断過程と、  
前記質問文情報から前記主たる名詞を認定する主名詞認定過程と、  
前記数量表現判断過程において質問文の解が数量表現であると判断した場合に、前記電子化テキスト情報から前記認定された主たる名詞と数量表現とが係り受け関係である組み合わせパターンを抽出し、抽出された前記組み合わせパターンから前記数量表現に付されている1又は複数の単位表現を抽出する単位表現抽出過程と、  
前記電子化テキスト情報における前記主たる名詞と前記抽出された各単位表現の出現頻度を算出し、算出した出現頻度情報に基づいて前記主たる名詞の一般的出現確率を算出する出現確率算出過程と、  
前記算出した主たる名詞の一般的出現確率に基づいて、前記電子化テキスト情報において前記各単位表現が出現するパターンのうち前記主たる名詞が前記各単位表現と共起して出現する回数の確率分布を算出する確率分布算出過程と、  
前記算出した確率分布における前記主たる名詞が前記各単位表現と共起して出現する回数が前記主たる名詞と数量表現とが係り受け関係である組み合わせパターンにおける前記数量表現に付されている各単位表現の出現回数以上である確率である検定確率を算出し、前記単位表現抽出過程において抽出した1又は複数の単位表現のうち、算出した前記検定確率が予め設定した閾値以上の単位表現を前記主たる名詞と共起して出現する可能性が低い不要な共起語と判断して除去し、前記検定確率が閾値未満の単位表現を有用な共起語と判断する不要単位表現除去過程と、  
前記有用と判断した単位表現を用いて前記質問文の解としての数量表現を推測する数量表現推測過程とを有することを特徴とする請求項1記載の統計的検定を利用した質問応答方法。

20

30

## 【請求項3】

前記表現推測過程は、  
前記質問文情報に基づいて、質問文の解が数量表現であるかを判断する数量表現判断過程と、  
前記質問文情報から前記主たる名詞を認定する主名詞認定過程と、  
前記数量表現判断過程において質問文の解が数量表現であると判断した場合に、前記電子化テキスト情報から前記認定された主たる名詞と数量表現とが係り受け関係である組み合わせパターンを抽出し、抽出された前記組み合わせパターンから前記数量表現に付されている1又は複数の単位表現を抽出する単位表現抽出過程と、  
前記電子化テキスト情報における前記主たる名詞と前記抽出された各単位表現の出現頻度を算出し、算出した出現頻度情報に基づいて前記各単位表現の一般的出現確率を算出する出現確率算出過程と、

40

50

前記各単位表現の一般的出現確率に基づいて、前記電子化テキスト情報において前記主たる名詞が出現するパターンのうち前記各単位表現が前記主たる名詞と共に起して出現する回数の確率分布を算出する確率分布算出過程と、  
前記算出した確率分布における前記各単位表現が前記主たる名詞と共に起して出現する回数が前記主たる名詞と数量表現とが係り受け関係である組み合わせパターンにおいて前記数量表現に付されている前記各単位表現と共に起して出現する前記数量表現の出現回数以上である確率である検定確率を算出し、前記単位表現抽出過程において抽出した1又は複数の単位表現のうち、算出した前記検定確率が予め設定した閾値以上の単位表現を前記主たる名詞と共に起して出現する可能性が低い不要な共起語と判断して除去し、前記検定確率が閾値未満の単位表現を有用な共起語と判断する不要単位表現除去過程と、  
前記有用と判断した単位表現を用いて前記質問文の解としての数量表現を推測する数量表現推測過程とを有する

10

ことを特徴とする請求項1記載の統計的検定を利用した質問応答方法。

【請求項4】

前記表現推測過程は、  
前記質問文情報に基づいて、質問文の解が数量表現であるかを判断する数量表現判断過程と、  
前記質問文情報から前記主たる名詞を認定する主名詞認定過程と、  
前記数量表現判断過程において質問文の解が数量表現であると判断した場合に、前記電子化テキスト情報から前記認定された主たる名詞と数量表現とが係り受け関係である組み合わせパターンを抽出し、抽出された前記組み合わせパターンから前記数量表現に付されている1又は複数の単位表現を抽出する単位表現抽出過程と、  
前記電子化テキスト情報における前記主たる名詞と前記抽出された各単位表現の出現頻度を算出し、算出した出現頻度情報に基づいて前記主たる名詞の一般的出現確率を算出する出現確率算出過程と、  
前記算出した主たる名詞の一般的出現確率に基づいて、前記電子化テキスト情報において前記各単位表現が出現するパターンのうち前記主たる名詞が前記各単位表現と共に起して出現する回数の確率分布を算出する確率分布算出過程と、  
前記算出した確率分布における前記主たる名詞が前記各単位表現と共に起して出現する回数が前記主たる名詞と数量表現とが係り受け関係である組み合わせパターンにおける前記数量表現に付されている各単位表現の出現回数以下である確率である検定確率を算出し、前記単位表現抽出過程において抽出した1又は複数の単位表現のうち、算出した前記検定確率が予め設定した閾値未満の単位表現を前記主たる名詞と共に起して出現する可能性が低い不要な共起語と判断して除去し、前記検定確率が閾値以上の単位表現を有用な共起語と判断する不要単位表現除去過程と、  
前記有用と判断した単位表現を用いて前記質問文の解としての数量表現を推測する数量表現推測過程とを有する

20

30

ことを特徴とする請求項1記載の統計的検定を利用した質問応答方法。

【請求項5】

前記表現推測過程は、  
前記質問文情報に基づいて、質問文の解が数量表現であるかを判断する数量表現判断過程と、  
前記質問文情報から前記主たる名詞を認定する主名詞認定過程と、  
前記数量表現判断過程において質問文の解が数量表現であると判断した場合に、前記電子化テキスト情報から前記認定された主たる名詞と数量表現とが係り受け関係である組み合わせパターンを抽出し、抽出された前記組み合わせパターンから前記数量表現に付されている1又は複数の単位表現を抽出する単位表現抽出過程と、  
前記電子化テキスト情報における前記主たる名詞と前記抽出された各単位表現の出現頻度を算出し、算出した出現頻度情報に基づいて前記各単位表現の一般的出現確率を算出する出現確率算出過程と、

40

50

前記各単位表現の一般的出現確率に基づいて、前記電子化テキスト情報において前記主たる名詞が出現するパターンのうち前記各単位表現が前記主たる名詞と共に起して出現する回数の確率分布を算出する確率分布算出過程と、  
 前記算出した確率分布における前記各単位表現が前記主たる名詞と共に起して出現する回数が前記主たる名詞と数量表現とが係り受け関係である組み合わせパターンにおいて前記数量表現に付されている前記各単位表現と共に起して出現する前記数量表現の出現回数以下である確率である検定確率を算出し、前記単位表現抽出過程において抽出した1又は複数の単位表現のうち、算出した前記検定確率が予め設定した閾値未満の単位表現を前記主たる名詞と共に起して出現する可能性が低い不要な共起語と判断して除去し、前記検定確率が閾値以上の単位表現を有用な共起語と判断する不要単位表現除去過程と、  
 前記有用と判断した単位表現を用いて前記質問文の解としての数量表現を推測する数量表現推測過程とを有する

10

ことを特徴とする請求項1記載の統計的検定を利用した質問応答方法。

【請求項6】

前記表現推測過程は、  
 前記質問文情報に基づいて、質問文の解が数量表現であるかを判断する数量表現判断過程と、  
 前記質問文情報から前記主たる名詞を認定する主名詞認定過程と、  
 前記数量表現判断過程において質問文の解が数量表現であると判断した場合に、前記電子化テキスト情報から前記認定された主たる名詞と数量表現とが係り受け関係である組み合わせパターンを抽出し、抽出された前記組み合わせパターンから前記数量表現に付されている1又は複数の単位表現を抽出する単位表現抽出過程と、  
 前記電子化テキスト情報における前記主たる名詞と前記抽出された各単位表現の出現頻度を算出する出現頻度抽出過程と、  
 超幾何分布を用いて、超幾何分布における前記主たる名詞が前記各単位表現と共に起して出現する回数が前記主たる名詞と数量表現とが係り受け関係である組み合わせパターンにおける前記数量表現に付されている各単位表現の出現回数以上である確率である検定確率を算出し、前記単位表現抽出過程において抽出した1又は複数の単位表現のうち、算出した前記検定確率が予め設定した閾値以上の単位表現を前記主たる名詞と共に起して出現する可能性が低い不要な共起語と判断して除去し、前記検定確率が閾値未満の単位表現を有用な共起語と判断する不要単位表現除去過程と、  
 前記有用と判断した単位表現を用いて前記質問文の解としての数量表現を推測する数量表現推測過程とを有する

20

30

ことを特徴とする請求項1記載の統計的検定を利用した質問応答方法。

【請求項7】

前記表現推測過程は、  
 前記質問文情報に基づいて、質問文の解が数量表現であるかを判断する数量表現判断過程と、  
 前記質問文情報から前記主たる名詞を認定する主名詞認定過程と、  
 前記数量表現判断過程において質問文の解が数量表現であると判断した場合に、前記電子化テキスト情報から前記認定された主たる名詞と数量表現とが係り受け関係である組み合わせパターンを抽出し、抽出された前記組み合わせパターンから前記数量表現に付されている1又は複数の単位表現を抽出する単位表現抽出過程と、  
 前記電子化テキスト情報における前記主たる名詞と前記抽出された各単位表現の出現頻度を算出する出現頻度抽出過程と、  
 超幾何分布を用いて、超幾何分布における前記主たる名詞が前記各単位表現と共に起して出現する回数が前記主たる名詞と数量表現とが係り受け関係である組み合わせパターンにおける前記数量表現に付されている各単位表現の出現回数以下である確率である検定確率を算出し、前記単位表現抽出過程において抽出した1又は複数の単位表現のうち、算出した前記検定確率が予め設定した閾値未満の単位表現を前記主たる名詞と共に起して出現する可

40

50

能性が低い不要な共起語と判断して除去し，前記検定確率が閾値以上の単位表現を有用な共起語と判断する不要単位表現除去過程と，  
前記有用と判断した単位表現を用いて前記質問文の解としての数量表現を推測する数量表現推測過程とを有することを特徴とする請求項 1 記載の統計的検定を利用した質問応答方法。

【請求項 8】

自然言語による質問文を入力し，データベースに記憶された電子化テキスト情報中の文との照合によって解を生成して出力する質問応答システムであって，  
前記質問文を入力する質問文入力手段と，  
入力された質問文の解が数量表現である場合に，統計的検定を用いて，前記電子化テキスト情報において出現する前記質問文の主たる名詞と数量表現との組み合わせパターンから抽出される単位表現のうち前記質問文中の主たる名詞と共起して出現する可能性が低い単位表現を不要な共起語と判断し，前記質問文中の主たる名詞と共起して出現する可能性が高い単位表現を前記質問文の解としての数量表現の推定に用いる有用な共起語と判断し，前記有用な共起語と判断した単位表現を用いて前記質問文の解としての数量表現を推測する表現推測手段と，  
前記推測された数量表現を用いて前記質問文の解を抽出する解答抽出手段とを備えることを特徴とする統計的検定を利用した質問応答システム。

10

【請求項 9】

前記表現推測手段は，  
前記質問文情報に基づいて，質問文の解が数量表現であるかを判断する数量表現判断手段と，  
前記質問文情報から前記主たる名詞を認定する主名詞認定手段と，  
前記数量表現判断手段が質問文の解が数量表現であると判断した場合に，前記電子化テキスト情報から前記認定された主たる名詞と数量表現とが係り受け関係である組み合わせパターンを抽出し，抽出された前記組み合わせパターンから前記数量表現に付されている 1 又は複数の単位表現を抽出する単位表現抽出手段と，  
前記電子化テキスト情報における前記主たる名詞と前記抽出された各単位表現の出現頻度を算出し，算出した出現頻度情報に基づいて前記主たる名詞の一般的出現確率を算出する出現確率算出手段と，  
前記算出した主たる名詞の一般的出現確率に基づいて，前記電子化テキスト情報において前記各単位表現が出現するパターンのうち前記主たる名詞が前記各単位表現と共起して出現する回数の確率分布を算出する確率分布算出手段と，  
前記算出した確率分布における前記主たる名詞が前記各単位表現と共起して出現する回数が前記主たる名詞と数量表現とが係り受け関係である組み合わせパターンにおける前記数量表現に付されている各単位表現の出現回数以上である確率である検定確率を算出し，前記単位表現抽出手段が抽出した 1 又は複数の単位表現のうち，算出した前記検定確率が予め設定した閾値以上の単位表現を前記主たる名詞と共起して出現する可能性が低い不要な共起語と判断して除去し，前記検定確率が閾値未満の単位表現を有用な共起語と判断する不要単位表現除去手段と，  
前記有用と判断した単位表現を用いて前記質問文の解としての数量表現を推測する数量表現推測手段とを備えることを特徴とする請求項 8 記載の統計的検定を利用した質問応答システム。

20

30

40

【請求項 10】

前記表現推測手段は，  
前記質問文情報に基づいて，質問文の解が数量表現であるかを判断する数量表現判断手段と，  
前記質問文情報から前記主たる名詞を認定する主名詞認定手段と，  
前記数量表現判断手段が質問文の解が数量表現であると判断した場合に，前記電子化テキスト情報から前記認定された主たる名詞と数量表現とが係り受け関係である組み合わせパ

50

ターンを抽出し、抽出された前記組み合わせパターンから前記数量表現に付されている 1 又は複数の単位表現を抽出する単位表現抽出手段と、  
前記電子化テキスト情報における前記主たる名詞と前記抽出された各単位表現の出現頻度を算出し、算出した出現頻度情報に基づいて前記各単位表現の一般的出現確率を算出する出現確率算出手段と、  
前記各単位表現の一般的出現確率に基づいて、前記電子化テキスト情報において前記主たる名詞が出現するパターンのうち前記各単位表現が前記主たる名詞と共に起して出現する回数の確率分布を算出する確率分布算出手段と、  
前記算出した確率分布における前記各単位表現が前記主たる名詞と共に起して出現する回数が前記主たる名詞と数量表現とが係り受け関係である組み合わせパターンにおいて前記数量表現に付されている前記各単位表現と共に起して出現する前記数量表現の出現回数以上である確率である検定確率を算出し、前記単位表現抽出手段が抽出した 1 又は複数の単位表現のうち、算出した前記検定確率が予め設定した閾値以上の単位表現を前記主たる名詞と共に起して出現する可能性が低い不要な共起語と判断して除去し、前記検定確率が閾値未満の単位表現を有用な共起語と判断する不要単位表現除去手段と、  
前記有用と判断した単位表現を用いて前記質問文の解としての数量表現を推測する数量表現推測手段とを備える  
ことを特徴とする請求項 8 記載の統計的検定を利用した質問応答システム。

10

**【請求項 11】**

前記表現推測手段は、  
前記質問文情報に基づいて、質問文の解が数量表現であるかを判断する数量表現判断手段と、  
前記質問文情報から前記主たる名詞を認定する主名詞認定手段と、  
前記数量表現判断手段が質問文の解が数量表現であると判断した場合に、前記電子化テキスト情報から前記認定された主たる名詞と数量表現とが係り受け関係である組み合わせパターンを抽出し、抽出された前記組み合わせパターンから前記数量表現に付されている 1 又は複数の単位表現を抽出する単位表現抽出手段と、  
前記電子化テキスト情報における前記主たる名詞と前記抽出された各単位表現の出現頻度を算出し、算出した出現頻度情報に基づいて前記主たる名詞の一般的出現確率を算出する出現確率算出手段と、  
前記算出した主たる名詞の一般的出現確率に基づいて、前記電子化テキスト情報において前記各単位表現が出現するパターンのうち前記主たる名詞が前記各単位表現と共に起して出現する回数の確率分布を算出する確率分布算出手段と、  
前記算出した確率分布における前記主たる名詞が前記各単位表現と共に起して出現する回数が前記主たる名詞と数量表現とが係り受け関係である組み合わせパターンにおける前記数量表現に付されている各単位表現の出現回数以下である確率である検定確率を算出し、前記単位表現抽出手段が抽出した 1 又は複数の単位表現のうち、算出した前記検定確率が予め設定した閾値未満の単位表現を前記主たる名詞と共に起して出現する可能性が低い不要な共起語と判断して除去し、前記検定確率が閾値以上の単位表現を有用な共起語と判断する不要単位表現除去手段と、  
前記有用と判断した単位表現を用いて前記質問文の解としての数量表現を推測する数量表現推測手段とを備える  
ことを特徴とする請求項 8 記載の統計的検定を利用した質問応答システム。

20

30

40

**【請求項 12】**

前記表現推測手段は、  
前記質問文情報に基づいて、質問文の解が数量表現であるかを判断する数量表現判断手段と、  
前記質問文情報から前記主たる名詞を認定する主名詞認定手段と、  
前記数量表現判断手段が質問文の解が数量表現であると判断した場合に、前記電子化テキスト情報から前記認定された主たる名詞と数量表現とが係り受け関係である組み合わせパ

50

ターンを抽出し、抽出された前記組み合わせパターンから前記数量表現に付されている 1 又は複数の単位表現を抽出する単位表現抽出手段と、  
 前記電子化テキスト情報における前記主たる名詞と前記抽出された各単位表現の出現頻度を算出し、算出した出現頻度情報に基づいて前記各単位表現の一般的出現確率を算出する出現確率算出手段と、  
 前記各単位表現の一般的出現確率に基づいて、前記電子化テキスト情報において前記主たる名詞が出現するパターンのうち前記各単位表現が前記主たる名詞と共に起して出現する回数の確率分布を算出する確率分布算出手段と、  
 前記算出した確率分布における前記各単位表現が前記主たる名詞と共に起して出現する回数が前記主たる名詞と数量表現とが係り受け関係である組み合わせパターンにおいて前記数量表現に付されている前記各単位表現と共に起して出現する前記数量表現の出現回数以下である確率である検定確率を算出し、前記単位表現抽出手段が抽出した 1 又は複数の単位表現のうち、算出した前記検定確率が予め設定した閾値未満の単位表現を前記主たる名詞と共に起して出現する可能性が低い不要な共起語と判断して除去し、前記検定確率が閾値以上の単位表現を有用な共起語と判断する不要単位表現除去手段と、  
 前記有用と判断した単位表現を用いて前記質問文の解としての数量表現を推測する数量表現推測手段とを備える  
 ことを特徴とする請求項 8 記載の統計的検定を利用した質問応答システム。

10

【請求項 13】

前記表現推測手段は、  
 前記質問文情報に基づいて、質問文の解が数量表現であるかを判断する数量表現判断手段と、  
 前記質問文情報から前記主たる名詞を認定する主名詞認定手段と、  
 前記数量表現判断手段が質問文の解が数量表現であると判断した場合に、前記電子化テキスト情報から前記認定された主たる名詞と数量表現とが係り受け関係である組み合わせパターンを抽出し、抽出された前記組み合わせパターンから前記数量表現に付されている 1 又は複数の単位表現を抽出する単位表現抽出手段と、  
 前記電子化テキスト情報における前記主たる名詞と前記抽出された各単位表現の出現頻度を算出する出現頻度抽出手段と、  
 超幾何分布を用いて、超幾何分布における前記主たる名詞が前記各単位表現と共に起して出現する回数が前記主たる名詞と数量表現とが係り受け関係である組み合わせパターンにおける前記数量表現に付されている各単位表現の出現回数以上である確率である検定確率を算出し、前記単位表現抽出手段が抽出した 1 又は複数の単位表現のうち、算出した前記検定確率が予め設定した閾値以上の単位表現を前記主たる名詞と共に起して出現する可能性が低い不要な共起語と判断して除去し、前記検定確率が閾値未満の単位表現を有用な共起語と判断する不要単位表現除去手段と、  
 前記有用と判断した単位表現を用いて前記質問文の解としての数量表現を推測する数量表現推測手段とを備える  
 ことを特徴とする請求項 8 記載の統計的検定を利用した質問応答システム。

20

30

【請求項 14】

前記表現推測手段は、  
 前記質問文情報に基づいて、質問文の解が数量表現であるかを判断する数量表現判断手段と、  
 前記質問文情報から前記主たる名詞を認定する主名詞認定手段と、  
 前記数量表現判断手段が質問文の解が数量表現であると判断した場合に、前記電子化テキスト情報から前記認定された主たる名詞と数量表現とが係り受け関係である組み合わせパターンを抽出し、抽出された前記組み合わせパターンから前記数量表現に付されている 1 又は複数の単位表現を抽出する単位表現抽出手段と、  
 前記電子化テキスト情報における前記主たる名詞と前記抽出された各単位表現の出現頻度を算出する出現頻度抽出手段と、

40

50

超幾何分布を用いて、超幾何分布における前記主たる名詞が前記各単位表現と共起して出現する回数が前記主たる名詞と数量表現とが係り受け関係である組み合わせパターンにおける前記数量表現に付されている各単位表現の出現回数以下である確率である検定確率を算出し、前記単位表現抽出手段が抽出した1又は複数の単位表現のうち、算出した前記検定確率が予め設定した閾値未満の単位表現を前記主たる名詞と共起して出現する可能性が低い不要な共起語と判断して除去し、前記検定確率が閾値以上の単位表現を有用な共起語と判断する不要単位表現除去手段と、  
前記有用と判断した単位表現を用いて前記質問文の解としての数量表現を推測する数量表現推測手段とを備える

ことを特徴とする請求項8記載の統計的検定を利用した質問応答システム。

10

【請求項15】

請求項1から請求項7までのいずれか1項に記載の統計的検定を利用した質問応答方法をコンピュータに実行させるための統計的検定を利用した質問応答プログラム。

【請求項16】

請求項1から請求項7までのいずれか1項に記載の統計的検定を利用した質問応答方法をコンピュータに実行させるための統計的検定を利用した質問応答プログラムを記録した記録媒体。

【発明の詳細な説明】

【0001】

【発明の属する技術分野】

本発明は、コンピュータによる自然言語の情報処理システムに係わり、特に統計的検定を利用して不要な共起語を除去する質問応答システムに関するものである。

20

【0002】

質問応答システムとは、例えば「日本の首都はどこですか」、「2002年のワールドカップの優勝国はどこですか」といった質問を入力すると、大量の電子化テキストから「東京」、「ブラジル」といった解を的確に出力するシステムのことである。

【0003】

質問応答システムは、検索した記事から解を探す必要がある情報検索などとは異なり、解自体を的確に出力するため、ユーザーがより早く解の情報を得ることができる。また、質問応答システムは、解自体を自動で出力するため、他の自動の知識処理システムの内部での知識処理システムとして利用することも可能である。本発明は、このように有用な質問応答システムのうち、統計的検定を利用して不要な共起語を除去する質問応答システムに関する。

30

【0004】

【従来の技術】

質問応答システムの実行例として、「20世紀最初にトラはどのくらいいましたか」という質問文を与えることを考える。当該質問文の解は「10万頭」であるが、これを求めるために一般の質問応答システムでは「どのくらい」という表現から数量表現が解であろうと推測する。また、さらに「トラはどのくらい」また「どのくらいいました」といった表現から、「どのくらい」に相当する数量表現としては、「...頭」といった「頭」を単位表現とした数量表現であろうと推測する。

40

【0005】

この情報に基づいて、例えば、キーワード抽出を行ない、「20世紀最初」「トラ」をキーワードとして記事や文書を検索し、その検索によって得られた記事や文書から、数字表現+「頭」のパターンを取り出すことで、解の「10万頭」を取り出すことができる。

【0006】

このような単位表現を含んだ解を取り出す質問応答システムは従来から存在した。(非特許文献1参照。)

【0007】

また、与えられた文書集合を特徴付ける単語を選出する方法についても、従来から存在し

50



ている（非特許文献2参照。）。

【0008】

【非特許文献1】

佐々木裕，磯崎秀樹，平博順，平尾努，賀沢秀人，鈴木潤，国領弘治，前田英作，S A I Q A：大量文書に基づく質問応答システム，情報処理学会自然言語処理研究会2001-NL-145，2001

【非特許文献2】

久光徹，丹羽芳樹，組み合わせ的確率モデルに基づく特徴単語選択方法，情報処理学会自然言語処理研究会，140-12，2000

【0009】

【発明が解決しようとする課題】

しかし，上記非特許文献1に記載された従来技術は，単位表現の抽出に人手で記述した規則，もしくはテーブルを利用しているため，質問文によっては解の抽出に有用な単位表現かを的確に判断することができず，低い正解率しか得られない場合が生じ得る。

【0010】

また，非特許文献2に記載された従来技術は，文章のキーワードとしての特徴単語の選択を目的としており，言語の機能的な表現である単位表現の抽出を意図したものではなかった。

【0011】

本発明は，上記従来技術の問題点を解決し，質問文に対する解が数量表現の場合に，当該数量表現に付される単位表現として適切な単位表現か否かを的確に判断することが可能な質問応答方法および質問応答システムを提供することを目的とする。

【0012】

【課題を解決するための手段】

上記課題を解決するため，本発明では，質問文中の主たる名詞と単位表現との新聞コーパス等の電子化テキスト情報における共起頻度情報に基づいて，統計的検定を用いることにより解の抽出に有用な単位表現かを判断し，当該有用な単位表現のみを解の抽出に用い，不要と判断した単位表現を解の抽出に用いないようにする。主たる名詞とは，質問文中の質問の主たる対象をいい，例えば上記「20世紀最初にトラはどのくらいいましたか」という質問文においては「トラ」が主たる名詞である。

【0013】

ここで，例えば「20世紀最初にトラはどのくらいいましたか」という質問文の「どのくらい」に相当する数量表現は，「．．．頭」といった「頭」を単位表現とした数量表現であろうと推測する場合，「トラは」+数量表現」といったパターンを大量に取り出し，その数量表現に付されている単位表現を抽出することで単位表現「頭」を取り出したり，「数量表現+「いました」」といったパターンから同様に単位表現「頭」を取り出したりすることができる。

【0014】

しかし，このような方法だけでは，例えば，「トラは1992年に．．．」といった文からは「年」という単位表現が取り出されることになる。「年」も単位表現として解の抽出に使えることとしてしまうと，「1992年」といった表現を誤って解と出力する可能性がある。

【0015】

また，コーパス等における共起頻度情報を用いる方法も考えられる。「トラ」と共起して出現する単位表現の回数を数え，この回数の最も大きな単位表現のみを用いるのである。このようにすると，おそらく「頭」がもっとも出現頻度が高いので，「頭」を単位表現とすることになり，「年」を単位表現として抽出してしまうという問題は解消される。

【0016】

しかし，「20世紀最初にトラはどのくらいいましたか」という質問文の「どのくらい」に相当する数量表現での単位表現としては，「匹」という表現も考えられる。解の表現が

10

20

30

40

50

「10万匹」となっていた場合，最大の頻度の「頭」だけを使うと，「10万匹」を解として取り出せなくなってしまう。

【0017】

そこで，本発明では，上記主たる名詞と単位表現とのコーパス等における共起頻度情報に基づき，統計的検定を用いて，前記数量表現に付されている各単位表現のうち前記主たる名詞と共起して出現する可能性が低い単位表現を不要な共起語と判断して除去するとともに，有用と判断した単位表現のみを解の抽出に用いるようにする。

【0018】

具体的には，

『当該単位表現と共起して出現する「トラは」の出現確率 = コーパスにおける「トラは」の一般的出現確率』 10

という仮説を立てて，コーパス等における「トラは」と単位表現とが共起して出現する回数や「トラは」と単位表現のそれぞれの出現頻度に基づいて上記仮説の検定を行うことを通じて，「トラは」と「頭」，「匹」等の各単位表現が偶然共起したもののか，必然的に共起しているもののかの判断を行なう。検定結果に基づいて，「トラは」とある単位表現が偶然共起したものであると判断される場合，それは偶然共起しただけであり，不要な単位表現と判断する。また，「トラは」とある単位表現が必然的に共起したものであると判断される場合，それは必然的に共起したものであるから，関係が深い表現であろうと予想されるので，有用な単位表現と判断する。この方法により，個々の単位表現に対して，解の抽出に用いる単位表現として有用か不要かの判断を下すことができる。そして，有用と判断された単位表現のみを用いて質問応答システムの解の判断を行なうと，先のすべての問題が解消されることになる。 20

【0019】

即ち，本発明は，自然言語による質問文を入力し，データベースに記憶された電子化テキスト情報中の文との照合によって解を生成して出力する質問応答方法であって，前記質問文を入力する質問文入力過程と，入力された質問文の解が数量表現である場合に，統計的検定を用いて，前記電子化テキスト情報において出現する前記質問文の主たる名詞と数量表現との組み合わせパターンから抽出される単位表現のうち前記質問文中の主たる名詞と共起して出現する可能性が低い単位表現を不要な共起語と判断し，前記質問文中の主たる名詞と共起して出現する可能性が高い単位表現を有用な共起語と判断し，前記有用な共起語と判断した単位表現を用いて前記質問文の解としての数量表現を推測する表現推測過程と，前記推測された数量表現を用いて前記質問文の解を抽出する解答抽出過程とを有することを特徴とする。 30

【0020】

また，本発明において，前記表現推測過程は，前記質問文情報に基づいて，質問文の解が数量表現であるかを判断する数量表現判断過程と，前記質問文情報から前記主たる名詞を認定する主名詞認定過程と，前記数量表現判断過程において質問文の解が数量表現であると判断した場合に，前記電子化テキスト情報から前記認定された主たる名詞と数量表現とが係り受け関係である組み合わせパターンを抽出し，抽出された前記組み合わせパターンから前記数量表現に付されている1又は複数の単位表現を抽出する単位表現抽出過程と，前記電子化テキスト情報における前記主たる名詞と前記抽出された各単位表現の出現頻度を算出し，算出した出現頻度情報に基づいて前記主たる名詞の一般的出現確率を算出する出現確率算出過程と，前記算出した主たる名詞の一般的出現確率に基づいて，前記電子化テキスト情報において前記各単位表現が出現するパターンのうち前記主たる名詞が前記各単位表現と共起して出現する回数の確率分布を算出する確率分布算出過程と，前記算出した確率分布における前記主たる名詞が前記各単位表現と共起して出現する回数が前記主たる名詞と数量表現とが係り受け関係である組み合わせパターンにおける前記数量表現に付されている各単位表現の出現回数以上である確率である検定確率を算出し，前記単位表現抽出過程において抽出した1又は複数の単位表現のうち，算出した前記検定確率が予め設定した閾値以上の単位表現を前記主たる名詞と共起して出現する可能性が低い不要な共起 40 50

語と判断して除去し、前記検定確率が閾値未満の単位表現を有用な共起語と判断する不要単位表現除去過程と、前記有用と判断した単位表現を用いて前記質問文の解としての数量表現を推測する数量表現推測過程とを有することを特徴とする。

【0021】

また、本発明において、前記表現推測過程は、前記質問文情報に基づいて、質問文の解が数量表現であるかを判断する数量表現判断過程と、前記質問文情報から前記主たる名詞を認定する主名詞認定過程と、前記数量表現判断過程において質問文の解が数量表現であると判断した場合に、前記電子化テキスト情報から前記認定された主たる名詞と数量表現とが係り受け関係である組み合わせパターンを抽出し、抽出された前記組み合わせパターンから前記数量表現に付されている1又は複数の単位表現を抽出する単位表現抽出過程と、前記電子化テキスト情報における前記主たる名詞と前記抽出された各単位表現の出現頻度を算出し、算出した出現頻度情報に基づいて前記各単位表現の一般的出現確率を算出する出現確率算出過程と、前記各単位表現の一般的出現確率に基づいて、前記電子化テキスト情報において前記主たる名詞が出現するパターンのうち前記各単位表現が前記主たる名詞と共起して出現する回数の確率分布を算出する確率分布算出過程と、前記算出した確率分布における前記各単位表現が前記主たる名詞と共起して出現する回数が前記主たる名詞と数量表現とが係り受け関係である組み合わせパターンにおいて前記数量表現に付されている前記各単位表現と共起して出現する前記数量表現の出現回数以上である確率である検定確率を算出し、前記単位表現抽出過程において抽出した1又は複数の単位表現のうち、算出した前記検定確率が予め設定した閾値以上の単位表現を前記主たる名詞と共起して出現する可能性が低い不要な共起語と判断して除去し、前記検定確率が閾値未満の単位表現を有用な共起語と判断する不要単位表現除去過程と、前記有用と判断した単位表現を用いて前記質問文の解としての数量表現を推測する数量表現推測過程とを有することを特徴とする。

【0022】

また、本発明において、前記表現推測過程は、前記質問文情報に基づいて、質問文の解が数量表現であるかを判断する数量表現判断過程と、前記質問文情報から前記主たる名詞を認定する主名詞認定過程と、前記数量表現判断過程において質問文の解が数量表現であると判断した場合に、前記電子化テキスト情報から前記認定された主たる名詞と数量表現とが係り受け関係である組み合わせパターンを抽出し、抽出された前記組み合わせパターンから前記数量表現に付されている1又は複数の単位表現を抽出する単位表現抽出過程と、前記電子化テキスト情報における前記主たる名詞と前記抽出された各単位表現の出現頻度を算出し、算出した出現頻度情報に基づいて前記主たる名詞の一般的出現確率を算出する出現確率算出過程と、前記算出した主たる名詞の一般的出現確率に基づいて、前記電子化テキスト情報において前記各単位表現が出現するパターンのうち前記主たる名詞が前記各単位表現と共起して出現する回数の確率分布を算出する確率分布算出過程と、前記算出した確率分布における前記主たる名詞が前記各単位表現と共起して出現する回数が前記主たる名詞と数量表現とが係り受け関係である組み合わせパターンにおける前記数量表現に付されている各単位表現の出現回数以下である確率である検定確率を算出し、前記単位表現抽出過程において抽出した1又は複数の単位表現のうち、算出した前記検定確率が予め設定した閾値未満の単位表現を前記主たる名詞と共起して出現する可能性が低い不要な共起語と判断して除去し、前記検定確率が閾値以上の単位表現を有用な共起語と判断する不要単位表現除去過程と、前記有用と判断した単位表現を用いて前記質問文の解としての数量表現を推測する数量表現推測過程とを有することを特徴とする。

【0023】

また、本発明において、前記表現推測過程は、前記質問文情報に基づいて、質問文の解が数量表現であるかを判断する数量表現判断過程と、前記質問文情報から前記主たる名詞を認定する主名詞認定過程と、前記数量表現判断過程において質問文の解が数量表現であると判断した場合に、前記電子化テキスト情報から前記認定された主たる名詞と数量表現とが係り受け関係である組み合わせパターンを抽出し、抽出された前記組み合わせパターン

から前記数量表現に付されている1又は複数の単位表現を抽出する単位表現抽出過程と、前記電子化テキスト情報における前記主たる名詞と前記抽出された各単位表現の出現頻度を算出し、算出した出現頻度情報に基づいて前記各単位表現の一般的出現確率を算出する出現確率算出過程と、前記各単位表現の一般的出現確率に基づいて、前記電子化テキスト情報において前記主たる名詞が出現するパターンのうち前記各単位表現が前記主たる名詞と共に出現する回数の確率分布を算出する確率分布算出過程と、前記算出した確率分布における前記各単位表現が前記主たる名詞と共に出現する回数が前記主たる名詞と数量表現とが係り受け関係である組み合わせパターンにおいて前記数量表現に付されている前記各単位表現と共に出現する前記数量表現の出現回数以下である確率である検定確率を算出し、前記単位表現抽出過程において抽出した1又は複数の単位表現のうち、算出した前記検定確率が予め設定した閾値未満の単位表現を前記主たる名詞と共に出現する可能性が低い不要な共起語と判断して除去し、前記検定確率が閾値以上の単位表現を有用な共起語と判断する不要単位表現除去過程と、前記有用と判断した単位表現を用いて前記質問文の解としての数量表現を推測する数量表現推測過程とを有することを特徴とする。

10

## 【0024】

また、本発明において、前記表現推測過程は、前記質問文情報に基づいて、質問文の解が数量表現であるかを判断する数量表現判断過程と、前記質問文情報から前記主たる名詞を認定する主名詞認定過程と、前記数量表現判断過程において質問文の解が数量表現であると判断した場合に、前記電子化テキスト情報から前記認定された主たる名詞と数量表現とが係り受け関係である組み合わせパターンを抽出し、抽出された前記組み合わせパターンから前記数量表現に付されている1又は複数の単位表現を抽出する単位表現抽出過程と、前記電子化テキスト情報における前記主たる名詞と前記抽出された各単位表現の出現頻度を算出する出現頻度抽出過程と、超幾何分布を用いて、超幾何分布における前記主たる名詞が前記各単位表現と共に出現する回数が前記主たる名詞と数量表現とが係り受け関係である組み合わせパターンにおける前記数量表現に付されている各単位表現の出現回数以上である確率である検定確率を算出し、前記単位表現抽出過程において抽出した1又は複数の単位表現のうち、算出した前記検定確率が予め設定した閾値以上の単位表現を前記主たる名詞と共に出現する可能性が低い不要な共起語と判断して除去し、前記検定確率が閾値未満の単位表現を有用な共起語と判断する不要単位表現除去過程と、前記有用と判断した単位表現を用いて前記質問文の解としての数量表現を推測する数量表現推測過程とを有することを特徴とする。

20

30

## 【0025】

また、本発明において、前記表現推測過程は、前記質問文情報に基づいて、質問文の解が数量表現であるかを判断する数量表現判断過程と、前記質問文情報から前記主たる名詞を認定する主名詞認定過程と、前記数量表現判断過程において質問文の解が数量表現であると判断した場合に、前記電子化テキスト情報から前記認定された主たる名詞と数量表現とが係り受け関係である組み合わせパターンを抽出し、抽出された前記組み合わせパターンから前記数量表現に付されている1又は複数の単位表現を抽出する単位表現抽出過程と、前記電子化テキスト情報における前記主たる名詞と前記抽出された各単位表現の出現頻度を算出する出現頻度抽出過程と、超幾何分布を用いて、超幾何分布における前記主たる名詞が前記各単位表現と共に出現する回数が前記主たる名詞と数量表現とが係り受け関係である組み合わせパターンにおける前記数量表現に付されている各単位表現の出現回数以下である確率である検定確率を算出し、前記単位表現抽出過程において抽出した1又は複数の単位表現のうち、算出した前記検定確率が予め設定した閾値未満の単位表現を前記主たる名詞と共に出現する可能性が低い不要な共起語と判断して除去し、前記検定確率が閾値以上の単位表現を有用な共起語と判断する不要単位表現除去過程と、前記有用と判断した単位表現を用いて前記質問文の解としての数量表現を推測する数量表現推測過程とを有することを特徴とする。

40

## 【0026】

50

また、本発明は、自然言語による質問文を入力し、データベースに記憶された電子化テキスト情報中の文との照合によって解を生成して出力する質問応答システムであって、前記質問文を入力する質問文入力手段と、入力された質問文の解が数量表現である場合に、統計的検定を用いて、前記電子化テキスト情報において出現する前記質問文の主たる名詞と数量表現との組み合わせパターンから抽出される単位表現のうち前記質問文中の主たる名詞と共起して出現する可能性が低い単位表現を不要な共起語と判断し、前記質問文中の主たる名詞と共起して出現する可能性が高い単位表現を前記質問文の解としての数量表現の推定に用いる有用な共起語と判断し、前記有用な共起語と判断した単位表現を用いて前記質問文の解としての数量表現を推測する表現推測手段と、前記推測された数量表現を用いて前記質問文の解を抽出する解答抽出手段とを備えることを特徴とする。

10

## 【0027】

また、本発明において、前記表現推測手段は、前記質問文情報に基づいて、質問文の解が数量表現であるかを判断する数量表現判断手段と、前記質問文情報から前記主たる名詞を認定する主名詞認定手段と、前記数量表現判断手段が質問文の解が数量表現であると判断した場合に、前記電子化テキスト情報から前記認定された主たる名詞と数量表現とが係り受け関係である組み合わせパターンを抽出し、抽出された前記組み合わせパターンから前記数量表現に付されている1又は複数の単位表現を抽出する単位表現抽出手段と、前記電子化テキスト情報における前記主たる名詞と前記抽出された各単位表現の出現頻度を算出し、算出した出現頻度情報に基づいて前記主たる名詞の一般的出現確率を算出する出現確率算出手段と、前記算出した主たる名詞の一般的出現確率に基づいて、前記電子化テキス

20

30

## 【0028】

また、本発明において、前記表現推測手段は、前記質問文情報に基づいて、質問文の解が数量表現であるかを判断する数量表現判断手段と、前記質問文情報から前記主たる名詞を認定する主名詞認定手段と、前記数量表現判断手段が質問文の解が数量表現であると判断した場合に、前記電子化テキスト情報から前記認定された主たる名詞と数量表現とが係り受け関係である組み合わせパターンを抽出し、抽出された前記組み合わせパターンから前記数量表現に付されている1又は複数の単位表現を抽出する単位表現抽出手段と、前記電子化テキスト情報における前記主たる名詞と前記抽出された各単位表現の出現頻度を算出し、算出した出現頻度情報に基づいて前記各単位表現の一般的出現確率を算出する出現確率算出手段と、前記各単位表現の一般的出現確率に基づいて、前記電子化テキスト情報に

40

50

## 【0029】

また、本発明において、前記表現推測手段は、前記質問文情報に基づいて、質問文の解が数量表現であるかを判断する数量表現判断手段と、前記質問文情報から前記主たる名詞を認定する主名詞認定手段と、前記数量表現判断手段が質問文の解が数量表現であると判断した場合に、前記電子化テキスト情報から前記認定された主たる名詞と数量表現とが係り受け関係である組み合わせパターンを抽出し、抽出された前記組み合わせパターンから前記数量表現に付されている1又は複数の単位表現を抽出する単位表現抽出手段と、前記電子化テキスト情報における前記主たる名詞と前記抽出された各単位表現の出現頻度を算出し、算出した出現頻度情報に基づいて前記主たる名詞の一般的出現確率を算出する出現確率算出手段と、前記算出した主たる名詞の一般的出現確率に基づいて、前記電子化テキスト情報において前記各単位表現が出現するパターンのうち前記主たる名詞が前記各単位表現と共に出現する回数の確率分布を算出する確率分布算出手段と、前記算出した確率分布における前記主たる名詞が前記各単位表現と共に出現する回数が前記主たる名詞と数量表現とが係り受け関係である組み合わせパターンにおける前記数量表現に付されている各単位表現の出現回数以下である確率である検定確率を算出し、前記単位表現抽出手段が抽出した1又は複数の単位表現のうち、算出した前記検定確率が予め設定した閾値未満の単位表現を前記主たる名詞と共に出現する可能性が低い不要な共起語と判断して除去し、前記検定確率が閾値以上の単位表現を有用な共起語と判断する不要単位表現除去手段と、前記有用と判断した単位表現を用いて前記質問文の解としての数量表現を推測する数量表現推測手段とを備えることを特徴とする。

10

20

## 【0030】

また、本発明において、前記表現推測手段は、前記質問文情報に基づいて、質問文の解が数量表現であるかを判断する数量表現判断手段と、前記質問文情報から前記主たる名詞を認定する主名詞認定手段と、前記数量表現判断手段が質問文の解が数量表現であると判断した場合に、前記電子化テキスト情報から前記認定された主たる名詞と数量表現とが係り受け関係である組み合わせパターンを抽出し、抽出された前記組み合わせパターンから前記数量表現に付されている1又は複数の単位表現を抽出する単位表現抽出手段と、前記電子化テキスト情報における前記主たる名詞と前記抽出された各単位表現の出現頻度を算出し、算出した出現頻度情報に基づいて前記各単位表現の一般的出現確率を算出する出現確率算出手段と、前記各単位表現の一般的出現確率に基づいて、前記電子化テキスト情報において前記主たる名詞が出現するパターンのうち前記各単位表現が前記主たる名詞と共に出現する回数の確率分布を算出する確率分布算出手段と、前記算出した確率分布における前記各単位表現が前記主たる名詞と共に出現する回数が前記主たる名詞と数量表現とが係り受け関係である組み合わせパターンにおいて前記数量表現に付されている前記各単位表現と共に出現する前記数量表現の出現回数以下である確率である検定確率を算出し、前記単位表現抽出手段が抽出した1又は複数の単位表現のうち、算出した前記検定確率が予め設定した閾値未満の単位表現を前記主たる名詞と共に出現する可能性が低い不要な共起語と判断して除去し、前記検定確率が閾値以上の単位表現を有用な共起語と判断する不要単位表現除去手段と、前記有用と判断した単位表現を用いて前記質問文の解としての数量表現を推測する数量表現推測手段とを備えることを特徴とする。

30

40

## 【0031】

また、本発明において、前記表現推測手段は、前記質問文情報に基づいて、質問文の解が数量表現であるかを判断する数量表現判断手段と、前記質問文情報から前記主たる名詞を認定する主名詞認定手段と、前記数量表現判断手段が質問文の解が数量表現であると判断した場合に、前記電子化テキスト情報から前記認定された主たる名詞と数量表現とが係り受け関係である組み合わせパターンを抽出し、抽出された前記組み合わせパターンから前記数量表現に付されている1又は複数の単位表現を抽出する単位表現抽出手段と、前記電子化テキスト情報における前記主たる名詞と前記抽出された各単位表現の出現頻度を算出する出現頻度抽出手段と、超幾何分布を用いて、超幾何分布における前記主たる名詞が前記各単位表現と共に出現する回数が前記主たる名詞と数量表現とが係り受け関係であ

50

る組み合わせパターンにおける前記数量表現に付されている各単位表現の出現回数以上である確率である検定確率を算出し、前記単位表現抽出手段が抽出した1又は複数の単位表現のうち、算出した前記検定確率が予め設定した閾値以上の単位表現を前記主たる名詞と共起して出現する可能性が低い不要な共起語と判断して除去し、前記検定確率が閾値未満の単位表現を有用な共起語と判断する不要単位表現除去手段と、前記有用と判断した単位表現を用いて前記質問文の解としての数量表現を推測する数量表現推測手段とを備えることを特徴とする。

#### 【0032】

また、本発明において、前記表現推測手段は、前記質問文情報に基づいて、質問文の解が数量表現であるかを判断する数量表現判断手段と、前記質問文情報から前記主たる名詞を認定する主名詞認定手段と、前記数量表現判断手段が質問文の解が数量表現であると判断した場合に、前記電子化テキスト情報から前記認定された主たる名詞と数量表現とが係り受け関係である組み合わせパターンを抽出し、抽出された前記組み合わせパターンから前記数量表現に付されている1又は複数の単位表現を抽出する単位表現抽出手段と、前記電子化テキスト情報における前記主たる名詞と前記抽出された各単位表現の出現頻度を算出する出現頻度抽出手段と、超幾何分布を用いて、超幾何分布における前記主たる名詞が前記各単位表現と共起して出現する回数が前記主たる名詞と数量表現とが係り受け関係である組み合わせパターンにおける前記数量表現に付されている各単位表現の出現回数以下である確率である検定確率を算出し、前記単位表現抽出手段が抽出した1又は複数の単位表現のうち、算出した前記検定確率が予め設定した閾値未満の単位表現を前記主たる名詞と共起して出現する可能性が低い不要な共起語と判断して除去し、前記検定確率が閾値以上の単位表現を有用な共起語と判断する不要単位表現除去手段と、前記有用と判断した単位表現を用いて前記質問文の解としての数量表現を推測する数量表現推測手段とを備えることを特徴とする。

10

20

#### 【0033】

また、本発明は、前記統計的検定を利用した質問応答方法をコンピュータに実行させるための統計的検定を利用した質問応答プログラムである。

#### 【0034】

また、本発明は、前記統計的検定を利用した質問応答方法をコンピュータに実行させるための統計的検定を利用した質問応答プログラムを記録した記録媒体である。

30

#### 【0035】

本発明を用いることにより、質問文に対する解が数量表現の場合に、当該数量表現に付される単位表現として適切な単位表現か否かを的確に判断することが可能な方法およびシステムを提供することが可能となる。

#### 【0036】

##### 【発明の実施の形態】

以下に、図を用いて、本発明の実施の形態を説明する。図1は本発明の統計的検定を利用した質問応答システムの構成の一例を示す図である。1は本発明の統計的検定を利用した質問応答システム、11は不要共起語除去システム、12は質問文情報が入力される質問文入力部、13はデータベース16からキーワード抽出や情報検索を行うキーワード抽出/情報検索部、14は質問文の解を抽出する解答抽出部、15は解答を出力する解答出力部、16は新聞コーパス等の電子化テキスト情報が記憶されたデータベース、17は質問文の解となり得そうな表現を推測する表現推測部、18はデータベース16から抽出した単位表現を記憶する単位表現データベースである。

40

#### 【0037】

図2は、不要共起語除去システム11の構成図の一例である。20は質問文入力部12に入力された質問文情報に基づいて、質問文の解が数量表現であるかを判断する数量表現判断手段、21は質問文情報から主たる名詞を認定する主名詞認定手段、22はデータベース16に記憶された電子化テキスト情報から主名詞認定手段21が認定した主たる名詞と数量表現とが係り受け関係である組み合わせパターンを抽出し、抽出された前記組み合わ

50

セパターンから前記数量表現に付されている 1 又は複数の単位表現を抽出し、単位表現データベース 18 に記録する単位表現抽出手段、23 は前記 1 又は複数の単位表現が前記主たる名詞と共に出現する可能性が高いか否かを判断する検定手段、24 は質問文の解としての数量表現を推測する数量表現推測手段である。

【0038】

また、230 は前記電子化テキスト情報における前記主たる名詞と前記抽出された各単位表現の出現頻度を算出し、算出した出現頻度情報に基づいて前記主たる名詞の一般的出現確率を算出する出現確率算出手段、231 は前記算出した主たる名詞の一般的出現確率に基づいて、前記電子化テキスト情報において前記各単位表現が出現するパターンのうち前記主たる名詞が前記各単位表現と共に出現する回数の確率分布を算出する確率分布算出手段、232 は確率分布算出手段 231 が算出した確率分布における前記主たる名詞が前記各単位表現と共に出現する回数が前記主たる名詞と数量表現とが係り受け関係である組み合わせパターンにおける前記数量表現に付されている各単位表現の出現回数以上である確率である検定確率を算出し、前記単位表現抽出手段が抽出した 1 又は複数の単位表現のうち、算出した前記検定確率が予め設定した閾値以上の単位表現を前記主たる名詞と共に出現する可能性が低い不要な共起語と判断して単位表現データベース 18 から除去し、前記検定確率が閾値未満の単位表現を有用な単位表現と判断する不要単位表現除去手段である。

10

【0039】

以下に、図 1、図 2 および図 8 を参照しつつ、図 3 および図 4 に基づいて本発明の第一の実施の形態を説明する。本発明の第一の実施の形態では、後述するように、各単位表現と共に出現する場合の主たる名詞のコーパスにおける出現確率  $p'$  が主たる名詞の一般的出現確率  $p$  と等しいという仮説を立てる。そして、一般的出現確率  $p$  に基づく主たる名詞がコーパスにおいて  $n$  回出現する各単位表現  $w_i$  と共に出現する回数が、実際にコーパスにおいて各単位表現  $w_i$  と主たる名詞とが共に出現する頻度である  $x_i$  回以上である確率（検定確率） $P$  の大きさに基づき上記仮説を右片側検定することを通じて  $p' > p$  かを結論付け、 $p' > p$  と結論付けできる単位表現は有用な単位表現と判断し、 $p' > p$  と結論付けできない単位表現は不要な単位表現と判断する。

20

【0040】

図 3 は、本発明の統計的検定を利用した質問応答処理フローの一例を示す図である。まず、質問文が質問文入力部 12 に入力される（ステップ S1）。例えば、「日本の国土面積はどのくらいですか」という質問文が質問文入力部 12 に入力される。次に、入力された質問文情報が質問文入力部 12 から不要共起語除去システム 11 の表現推測部 17 に渡される（ステップ S2）。次に、不要共起語除去システム 11 において、表現推測部 17 が、統計的検定を用いて、コーパスから抽出した単位表現のうち、解となり得そうな表現の推測に有用な単位表現かまたは不要な単位表現かを判断し、不要と判断した単位表現を除去する（ステップ S3）。

30

【0041】

次に、表現推測部 17 が、有用と判断された単位表現を用いて、解となり得そうな表現を推測する（ステップ S4）。具体的には、表現推測部 17 の数量表現推測手段 24 が、解となり得そうな数量表現を推測する。

40

【0042】

そして、キーワード抽出/情報検索部 13 が、表現推測部 17 から渡された質問文からキーワードを抽出する。そして、抽出したキーワードを用いてデータベースから解が記述してありそうな記事群を取り出し、取り出した記事群を解答抽出部 14 に渡す（ステップ S5）。

【0043】

次に、解答抽出部 14 において、前記取り出された記事群から、表現推測部 17 で推測した表現に合致する表現を抽出し、抽出した表現を解答出力部 15 に渡す（ステップ S6）。最後に、渡された表現を解答出力部 15 が解答として出力する（ステップ S7）。

50



## 【0044】

ここで、本発明の第一の実施の形態における不要な単位表現の除去処理フローの詳細の一例を図4に示す。図4は、図3のステップS3の詳細を示す図である。まず、数量表現判断手段20が、前記入力された質問文情報に基づいて、質問文の解が数量表現であるかを判断する(ステップS21)。例えば、質問文が「日本の国土面積はどのくらいですか」の場合のように「どのくらい」などの表現を含んでいた場合は、質問文の解が数量表現であると判断される。質問文の解が数量表現でないと判断する場合は前記ステップS5へ進む。

## 【0045】

質問文の解が数量表現であると判断する場合は、数量表現判断手段20は、質問文が単位表現を有していないかを判断する(ステップS22)。質問文が単位表現を有していない場合には、主名詞認定手段21は、質問文情報から質問文における主たる名詞を認定する(ステップS23)。主たる名詞の認定は、「Xはどのくらい」のパターンから、形態素解析と文パターンを用いた規則に基づいてなされる。例えば上記質問文における主たる名詞は、「面積」として認定される。質問文が単位表現を有している場合には、処理を終了する。

## 【0046】

次に、検定手段23の出現確率算出手段230は、前記コーパスにおける前記主たる名詞の一般的出現確率 $p$ を算出する(ステップS24)。かかる一般的出現確率 $p$ は、コーパスにおける前記主たる名詞の出現回数をコーパスの規模で割ったものである。毎日新聞などの大規模な新聞コーパスを用いると、例えば、上記「面積」のコーパスにおける出現回数は7,933回となる。コーパスの規模が、例えば409,502,077文字であるとすると、「面積」の一般的出現確率 $p$ は、

$$p = 7,933 / 409,502,077$$

$$= 0.00001937230711530676$$

と算出される。

## 【0047】

次に、単位表現抽出手段22が、コーパスから前記主たる名詞と数量表現とが係り受け関係である組み合わせパターンを抽出し、抽出された前記組み合わせパターンから前記数量表現に付されている1又は複数の単位表現 $w_i$ を抽出し、単位表現データベース18に記録する(ステップS25)。本実施の形態では、「面積は」+数量表現のパターンをコーパスから取り出し、数量表現に付されている単位表現 $w_i$ を取り出す。各単位表現 $w_i$ の取り出しは、形態素解析結果、また、単語辞書を用いて行なわれる。抽出された各単位表現 $w_i$ は、単位表現データベース18に記憶される。数量表現のパターンの取り出しは形態素解析と文パターンを用いた規則に基づいてなされる。

## 【0048】

ここで、コーパスから「面積は」+数量表現のパターンで取り出した数量表現から以下の単位表現 $w_i$ が例えば以下の頻度で抽出できる。

## 【0049】

平方メートル：41回、ヘクタール：27回、平方キロメートル：5回、倍：2回、分：2回、畳：1回、番：1回、平方センチメートル：1回  
 コロンの前が単位表現 $w_i$ で、数字が頻度である。「面積は」+数量表現のパターンの頻度は、上記の頻度の合計で80回である。

## 【0050】

次に、出現確率算出手段230が、上記抽出した各単位表現 $w_i$ のコーパスにおける出現頻度 $n$ を算出する(ステップS26)。上記各単位表現 $w_i$ のコーパスでの出現頻度 $n$ は、例えば平方メートル：17,510回、ヘクタール：8,088回、平方キロメートル：247回、倍：41,686回、分：730,790回、畳：4,829回、番：124,233回、平方センチメートル：20回である。

## 【0051】

また、出現確率算出手段 230 は、実際にコーパスにおいて各単位表現  $w_i$  と主たる名詞とが共起して出現する頻度  $x_i$  を算出する（ステップ S27）。ここでの共起の定義は、例えば主たる名詞が「面積」の場合、「面積は」+数量表現のパターンにおける数量表現の単位表現が  $w_i$  であることとする。即ち、コーパスにおける各単位表現  $w_i$  の  $n$  個の出現パターンのうち、その数量表現と「面積」が「面積は」+数量表現のパターンで共起して出現した頻度が  $x_i$  である。例えば、平方メートルの例だと、 $x_i = 41$  回となる。

#### 【0052】

次に、「平方メートル」と共起する場合の「面積」の出現確率  $p'$  = そのような条件のないときの「面積」の一般的出現確率  $p$  という仮説を立てて、二項分布を利用した検定を行う。即ち、まず、確率分布算出手段 231 が、前記算出した主たる名詞の一般的出現確率  $p$  に基づく前記コーパスにおいて前記主たる名詞と前記単位表現  $w_i$  とが共起して出現する回数の確率分布を算出する（ステップ S28）。個々の試行はすべて独立とし、ある単位表現  $w_i$  が  $n$  回出現しているときに、一般的出現確率  $p$  の主たる名詞が出現する回数の確率分布を求めるのである。かかる確率分布は、主たる名詞の出現回数を  $r$  ( $r = 0, 1, 2, 3, \dots, n$ ) とし、

$${}_n C_r p^r (1-p)^{n-r}$$

と算出され、図 8 に示されるような二項分布となる。

#### 【0053】

次に、不要単位表現除去手段 232 は、算出した確率分布における前記主たる名詞と各単位表現  $w_i$  とが共起して出現する回数が  $x_i$  回以上である確率（検定確率） $P$  を算出する（ステップ S29）。図 8 では斜線部分の確率が  $P$  である。例えば、「平方メートル」の  $n = 17510$  回の出現場面において、一般的出現確率  $p = \text{約} 0.00001937$  の「面積」が、 $x_i = 41$  回以上現れる検定確率  $P$  を求める。

#### 【0054】

本発明の実施の形態では、「平方メートル」と共起する場合の「面積」の出現確率  $p'$  とそのような条件のないときの「面積」の一般的出現確率  $p$  は本来異なる可能性があるのに同じであると仮定しており、「平方メートル」と共起する場合の「面積」の出現確率  $p'$  に、そのような条件のないときの「面積」の一般的出現確率  $p$  を用いている。

#### 【0055】

従って、検定確率  $P$  の値が十分に小さい場合は、 $p' = p$  という仮説は棄却され、 $p'$  と  $p$  とは異なると判断できる。また、本発明の実施の形態では、「 $x_i$  回以上」のように、検定確率  $P$  を求める際に、確率分布の片側の領域のみを用いているので、片側検定になっている。さらに、用いている領域が「 $x_i$  回以上」のように大きい場合の方の領域を用いる右片側検定であるため、 $P$  の値が十分に小さい場合、 $p' > p$  と結論付けできる。ただし、かかる判断は必ず正しいという意味ではなく、確率  $P$  だけは誤る可能性をもった判断である。

#### 【0056】

不要単位表現除去手段 232 は、 $P$  の値が十分に小さく、「平方メートル」と共起する場合の「面積」の出現確率  $p'$  が、そのような条件のないときの「面積」の一般的出現確率  $p$  よりも大きいと結論付けできる場合は、この「平方メートル」と「面積」の共起は偶然ではなく必然的な共起であるとして「平方メートル」を有用な単位表現と判断する。

#### 【0057】

逆に、検定確率  $P$  の値が十分に小さくなく、「平方メートル」と共起する場合の「面積」の出現確率  $p'$  が、そのような条件のないときの「面積」の一般的出現確率  $p$  よりも大きいと結論できない場合は、この「平方メートル」と「面積」の共起は必然的なものではなく偶然的な共起である可能性が高いとして「平方メートル」を不要な単位表現と判断し、以降の質問応答システムの処理では用いないようにする。

#### 【0058】

即ち、不要単位表現除去手段 232 は、各単位表現  $w_i$  の検定確率  $P$  の値が予め設定し

10

20

30

40

50

た閾値未満かを判断し(ステップS30), Pの値が閾値未満の単位表現を有用な単位表現と判断し(ステップS32), Pの値が閾値以上の単位表現を不要な単位表現と判断して単位表現データベース18から除去し(ステップS31), 処理を終了する。

【0059】

実際に上記の例でPの値を計算すると, 平方メートル: 0.000000000, ヘクタール: 0.000000000, 平方キロメートル: 0.000000000, 倍: 0.19392510, 分: 0.99998923, 畳: 0.08930732, 番: 0.90988811, 平方センチメートル: 0.00038737となる。この計算は上述したように, 二項分布の理論を使うことで計算することができる。ここでは, フリーソフトのMath-CDF-0.1のサブルーチンpbinomを用いて計算した。

10

【0060】

ここで, Pの値が十分小さいことを意味する閾値として, 例えば, 0.1を用いると,

平方メートル: 0.000000000

ヘクタール: 0.000000000

平方キロメートル: 0.000000000

畳: 0.08930732

平方センチメートル: 0.00038737

が有用な単位表現, すなわち妥当な共起語で,

倍: 0.19392510

分: 0.99998923

番: 0.90988811

が不要な単位表現, すなわち妥当でない共起語であると判断できる。

20

【0061】

実際, 平方メートル: 0.000000000, ヘクタール: 0.000000000, 平方キロメートル: 0.000000000, 畳: 0.08930732, 平方センチメートル: 0.00038737は, 面積の単位表現として利用できるが, 倍: 0.19392510, 分: 0.99998923, 番: 0.90988811は, 面積の単位表現としては不当な表現である。

【0062】

このようにして妥当な単位表現と判断した単位表現を用いて質問文の解としての数量表現の推測に利用するのである。

30

【0063】

本発明の第一の実施の形態においては, 「平方メートル」と共起する場合の「面積」の出現確率 $p'$ が「面積」の一般的出現確率 $p$ と等しいという仮説を立てて, Pの値の大きさに基づいて右片側検定を行って, 上記仮説が棄却できるか, 即ち $p' > p$ かを結論付け,  $p' > p$ と結論付けできる単位表現は有用な単位表現であり,  $p' > p$ と結論付けできない単位表現は不要な単位表現と判断していたが, 本発明においては, 以下のように左片側検定を行って $p' < p$ かを結論付け,  $p' < p$ と結論付けできる単位表現は不要な単位表現であり,  $p' < p$ と結論付けできない単位表現は有用な単位表現と判断することもできる。

40

【0064】

本発明の第二の実施の形態においては, まず第一の実施の形態と同様に「面積」の表現の一般的出現確率 $p$ と, 各単位表現の出現数 $n$ を求める。「平方メートル」の例では例えば $n = 17, 510$ となる。

【0065】

この $n$ 個のパターンのうち, その数量表現と「面積」が「面積は」+数量表現のパターンで共起して出現した頻度を $x_i$ とする。例えば, 平方メートルの例だと,  $x_i = 41$ となる。

【0066】

次に, 個々の試行はすべて独立と仮定し, 「平方メートル」の $n$ 回の出現場面において,

50

一般的出現確率  $p$  の「面積」が、 $x_i$  回以下現れる確率（検定確率） $P$  を求める。言い換えると、「平方メートル」が  $n$  回出現していて、1 回に  $p$  の確率で出現する「面積」が、この  $n$  回の「平方メートル」とともに共起して出現する回数が  $x_i$  回以下である確率（検定確率） $P$  を求める。

【0067】

ここでは、「平方メートル」と共起する場合の「面積」の出現確率  $p'$  とそのような条件のないときの「面積」の一般的出現確率  $p$  は本来異なる可能性があるのに同じであると仮定しており、「平方メートル」と共起する場合の「面積」の出現確率  $p'$  に、そのような条件のないときの「面積」の一般的出現確率  $p$  を用いている。

【0068】

従って、検定確率  $P$  の値が十分に小さい場合は、 $p' = p$  という仮説は棄却され、 $p'$  と  $p$  とは異なると判断できる。また、本発明の実施の形態では、「 $x_i$  回以下」のように、検定確率  $P$  を求める際に、確率分布の片側の領域のみを用いているので、片側検定になっている。さらに、用いている領域が「 $x_i$  回以下」のように小さい場合の方の領域を用いる左片側検定であるので、 $P$  の値が十分に小さい場合、 $p' < p$  と結論付けできる。ただし、かかる判断は必ず正しいという意味ではなく、確率  $P$  だけは誤る可能性をもった判断である。

【0069】

不要単位表現除去手段 232 は、 $P$  の値が十分に小さく、「平方メートル」と共起する場合の「面積」の出現確率  $p'$  が、そのような条件のないときの「面積」の一般的出現確率  $p$  よりも小さいと結論付けできる場合は、この「平方メートル」と「面積」の共起は偶然であるとして「平方メートル」を不要な単位表現と判断し、単位表現データベース 18 から除去する。

【0070】

逆に、検定確率  $P$  の値が十分に小さくなく、「平方メートル」と共起する場合の「面積」の出現確率  $p'$  が、そのような条件のないときの「面積」の一般的出現確率  $p$  よりも小さいと結論付けできない場合は、この「平方メートル」と「面積」の共起は偶然的なものではなく必然的な共起である可能性が高いとして「平方メートル」を有用な単位表現と判断し、当該有用と判断した単位表現を用いて数量表現推測手段 24 が質問文の解としての数量表現を推測する。

【0071】

以上整理すると、 $P$  の値が小さいほど、不要な単位表現と判断し、 $P$  の値が大きいほど、有用な単位表現と判断するということになる。

【0072】

実際に上記の例で  $P$  の値を計算すると、

平方メートル：1.000000000

ヘクタール：1.000000000

平方キロメートル：1.000000000

倍：0.95148679

分：0.00008198

畳：0.99588861

番：0.30698656

平方センチメートル：0.99999993

となる。

【0073】

この計算は二項分布の理論を使うことで計算することができる。ここでは、フリーソフトの Math-CDF-0.1 のサブルーチン `pbinom` を用いて計算した。

【0074】

ここで、 $P$  の値が十分小さいことを意味する閾値として、例えば、0.99 を用いると、 $P$  の値が 0.99 未満の単位表現は不要な単位表現であって不要な共起語であると判断し

10

20

30

40

50

、Pの値が0.99以上の単位表現は主たる名詞「面積」と共起して出現する可能性が高く、有用な単位表現であって妥当な共起語であると判断する。即ち、

平方メートル：1.000000000

ヘクタール：1.000000000

平方キロメートル：1.000000000

畳：0.99588861

平方センチメートル：0.9999993

が妥当な共起語で、

倍：0.95148679

分：0.00008198

番：0.30698656

が妥当でない不要な共起語であると判断される。

【0075】

実際、平方メートル：1.000000000、ヘクタール：1.000000000、平方キロメートル：1.000000000、畳：0.99588861、平方センチメートル：0.9999993は面積の単位表現として利用できるが、倍：0.95148679、分：0.00008198、番：0.30698656は、面積の単位表現としては不当な表現である。

【0076】

このようにして求めた妥当な単位表現のみを用いて質問文の解としての数量表現の推測に利用するのである。

【0077】

図5は、上記本発明の第二の実施の形態における不要な単位表現の除去処理フローの詳細の一例を示す図であり、図3のステップS3の詳細であるステップS41乃至ステップS52を示したものである。なお、本発明の第二の実施の形態においては、図3におけるステップ3以外のステップは第一の実施の形態と同様である。

【0078】

図5に示すように、本発明の第二の実施の形態においては、ステップS49において、算出した確率分布における主たる名詞と各単位表現 $w_i$ とが共起して出現する回数が $x_i$ 回以下である確率P(検定確率)を算出することと、ステップS52においてPの値< 30  
閾値の場合には不要な単位表現と判断して単位表現データベースから除去することと、ステップS51においてPの値<閾値でない場合には有用な単位表現と判断する点において図4に示す本発明の第一の実施の形態と異なり、図5における他のステップは本発明の第一の実施の形態と同様である。

【0079】

本発明は、その趣旨に基づき、以下のように種々の変形が可能である。

【0080】

上記本発明の第一の実施の形態および第二の実施の形態では、一般的出現確率をpとして、コーパスでの主たる名詞としての「面積」の出現確率を用いたが、本発明では、コーパスにおける一般的出現確率pとしては、単位表現 $w_i$ の出現確率を用いて、コーパスにおいて主たる名詞「面積」が出現した個数をnとしても同様の検定が行なえる。

【0081】

即ち、「面積」と共起する場合の「平方メートル」の出現確率 $p'$ が「平方メートル」の一般的出現確率pと等しいという仮説を立てて左片側検定を行い、「平方メートル」が上記n回出現する主たる名詞「面積」と共起して出現する回数が $x_i$ 回(例えば41回)以下である検定確率Pの値の大きさが閾値(例えば0.99)未満である場合は $p' < p$ であると結論付けて「平方メートル」は不要な単位表現と判断し、検定確率Pの値の大きさが閾値以上であり、 $p' < p$ と結論付けできない場合は有用な単位表現と判断することもできる。

【0082】

10

20

30

40

50

例えば、「平方メートル」の場合だと、

$p = 17510 / 409502077$ ,  $n = 7933$ であり、

$P = 1.000000000 > 0.99$

となつて、 $p' < p$ と判断できないため、「平方メートル」は有用な単位表現であると判断できる。

【0083】

かかる左片側検定を行う本発明の第三の実施の形態の不要な単位表現の除去処理フローの詳細の一例を図6に示す。図6は、図3のステップS3の詳細であるステップS81乃至ステップS92を示したものである。図6に示すように、本発明の第三の実施の形態においては、ステップS84においてコーパスにおける各単位表現 $w_i$ の一般的出現確率 $p$ を算出していることと、ステップS86において、主たる名詞のコーパスにおける出現頻度 $n$ を算出している点において図5に示す本発明の第二の実施の形態と異なる。

10

【0084】

もちろん、本発明においては、「面積」と共起する場合の「平方メートル」の出現確率 $p'$ が「平方メートル」の一般的出現確率 $p$ と等しいという仮説を立てて右片側検定を行い、「平方メートル」が上記 $n$ 回出現する「面積」と共起して出現する回数が $x_i$ 回以上である検定確率 $P$ の値の大きさが閾値未満である場合は $p' > p$ と結論付けて有用な単位表現と判断し、検定確率 $P$ の値の大きさが閾値以上であり、 $p' > p$ と結論付けできない場合は不要な単位表現と判断することもできる。

【0085】

かかる右片側検定を行う本発明の第四の実施の形態の不要な単位表現の除去処理フローの詳細の一例を図7に示す。図7は、図3のステップS3の詳細であるステップS61乃至ステップS72を示したものである。図7に示すように、本発明の第四の実施の形態においては、ステップS64においてコーパスにおける各単位表現 $w_i$ の一般的出現確率 $p$ を算出していることと、ステップS66において、主たる名詞のコーパスにおける出現頻度 $n$ を算出している点において図4に示す本発明の第一の実施の形態と異なる。

20

【0086】

また、本発明においては、以下に示す超幾何分布を用いた検定を行うこともできる。

【0087】

超幾何分布とは、

$$hg(N, k, n, m) = C(K, m) \times C(N - K, n - m) / C(N, n)$$

の形で表せる分布である。

30

【0088】

ただし、 $C$ は「組み合わせ」を意味する記号 $C(A, B) = A! / B! / (A - B)!$ であり、また、 $hg(N, k, n, m)$ は、「 $N$ 個の玉の中に $K$ 個の赤い玉があるとき、任意に取り出した $n$ 個の玉の中に赤い玉がちょうど $m$ 個含まれる確率」である。

【0089】

ここで、 $hgs(N, K, n, k) = hg(N, K, n, m)$

とすると、 $hgs(N, K, n, k)$ は、「 $N$ 個の玉の中に $K$ 個の赤い玉があるとき、任意に取り出した $n$ 個の玉の中に赤い玉が $k$ 個以上含まれる確率」となる。なお、 $hg(N, K, n, m)$ は、 $m \geq k$ である $m$ についての $hg(N, k, n, m)$ の合計を意味する。

40

【0090】

ここで、 $N, K, n, k$ を以下のように解釈する。

【0091】

$N$ ：コーパスの大きさ

$K$ ：「面積」の出現頻度

$n$ ：ある単位表現の出現頻度

$k$ ：「面積」と、ある単位表現の共起回数

上記解釈によると、 $hgs(N, K, n, k)$ は、「 $N$ の大きさのコーパスの中に「面積

50

」という表現が  $K$  個あるときに、ある着目している単位表現を  $n$  個取り出し、その  $n$  個の単位表現のうち、 $k$  個以上のもので、「面積」とその単位表現が共起する事象の起る確率」を意味することとなる。

【0092】

これは、前記本発明の実施の形態において、 $p = K / N$  を「面積」の一般的出現確率としていたところを、 $K / N$  の形にまとめずに  $K$  と  $N$  にわけたまま扱っていることに相当する。超幾何分布を用いる方法では、 $p = K / N$  とせずに  $K$  と  $N$  とにわけている分、仮定が少なく、近似の少ない手法で精度は高くなるものである。

【0093】

この方法では、 $hgs(N, K, n, k)$  が小さいほど、「面積」と単位表現の共起が妥当なものとして判断でき、 $hgs(N, K, n, k)$  が大きいほど、「面積」と単位表現の共起は偶然のもので、妥当なものではないと判断できる。

10

【0094】

従って、本発明においては、不要単位表現除去手段 232 は、 $hgs(N, K, n, k)$  がある閾値より小さい単位表現を有用な単位表現と判断して質問文の解としての数量表現の推定に用いるようにし、 $hgs(N, K, n, k)$  がある閾値以上の単位表現を不要な単位表現と判断して除去する。

【0095】

本発明の第五の実施の形態では、上述した超幾何分布を用いて右片側検定を行う。図 9 は、本発明の第五の実施の形態における不要な単位表現の除去処理フローの詳細の一例を示す図であり、図 3 のステップ S3 の詳細であるステップ S101 乃至ステップ S111 を示したものである。なお、図 3 におけるその他のステップは本発明の第一の実施の形態と同様である。

20

【0096】

図 9 に示すように、本発明の第五の実施の形態においては、質問文の解が数量表現であるかを判断し（ステップ S101）、解が数量表現でない場合は図 3 のステップ S5 に移行し、解が数量表現の場合は質問文が単位表現を有していないかを判断する（ステップ S102）。

【0097】

質問文が単位表現を有している場合には処理を終了して図 3 のステップ S4 に移行し、質問文が単位表現を有していない場合には、質問文情報から質問文における主たる名詞を認定する（ステップ S103）。

30

【0098】

次に、コーパスにおける主たる名詞の出現頻度  $K$  を算出する（ステップ S104）。そして、コーパスから主たる名詞と数量表現とが係り受け関係である組み合わせパターンを抽出し、抽出された前記組み合わせパターンから前記数量表現に付されている 1 又は複数の単位表現  $w_i$  を抽出し、単位表現データベース 18 に記録する（ステップ S105）。そして、抽出した各単位表現  $w_i$  のコーパスにおける出現頻度  $n$  を算出する（ステップ S106）。

【0099】

次に、各単位表現  $w_i$  と主たる名詞とが共起して出現する回数  $k$  を算出する（ステップ S107）。次に、超幾何分布における主たる名詞と各単位表現  $w_i$  とが共起して出現する回数が  $k$  回以上である検定確率  $P = hgs(N, K, n, k)$  を算出する（ステップ S108）。なお、 $N$  はコーパスの大きさである。

40

【0100】

そして、 $P$  の値 < 閾値であるかを判断し（ステップ S109）、 $P$  の値 < 閾値である場合には有用な単位表現と判断し（ステップ S111）、 $P$  の値 < 閾値でない場合は、不要な単位表現と判断して単位表現データベースから除去する（ステップ S110）。

【0101】

また、本発明においては、超幾何分布を用いて左片側検定を行うこともできる。超幾何分

50

布を用いて左片側検定を行う本発明の第六の実施の形態における不要単位表現の除去処理フローの詳細の一例を図10に示す。図10は、図3のステップS3の詳細であるステップS121乃至ステップS131を示したものである。本発明の第六の実施の形態においては、図10のステップS128において超幾何分布における主たる名詞と各単位表現  $w_i$  とが共起して出現する回数が  $k$  回以下である検定確率  $P = h g s ( N , K , n , k )$  を算出することと、ステップS130において、 $P$  の値  $<$  閾値でない場合には有用な単位表現と判断することと、ステップS131において、 $P$  の値  $<$  閾値である場合は、不要な単位表現と判断して単位表現データベースから除去する点で図9に示す本発明の第五の実施の形態と異なる。

#### 【0102】

なお、本発明の第五の実施の形態または本発明の第六の実施の形態においては、 $K =$  主たる名詞の出現頻度、 $n =$  各単位表現の出現頻度としたが、この  $K$  と  $n$  は交換可能で、 $K =$  各単位表現の出現頻度、 $n =$  主たる名詞の出現回数としてもよい。

#### 【0103】

以上説明した全ての本発明の実施の形態では、ある種の比率の検定を行なっていることに相当する。おおよそ、主たる名詞と単位表現の一般的出現確率から予想される主たる名詞と単位表現の共起回数よりも大きい回数もしくは共起回数以上の回数で主たる名詞と単位表現の共起が出現しているかどうかを検定するのである。

#### 【0104】

また、本発明は、AIC や Z - s c o r e などの比率の検定を行なうことができる統計的

#### 【0105】

なお、本発明は、前記従来技術における人手で記述した規則、もしくは、テーブルの作成の補助にも用いることができる。

#### 【0106】

#### 【発明の効果】

本発明を用いることにより、質問文に対する解が数量表現の場合に、当該数量表現に付される単位表現として適切な単位表現か否かを的確に判断することが可能な方法およびシステムを提供することが可能となる。

#### 【図面の簡単な説明】

【図1】統計的検定を利用した質問応答システムの構成の一例を示す図である。

【図2】不要共起語除去システムの構成図の一例である。

【図3】統計的検定を利用した質問応答処理フローの一例を示す図である。

【図4】不要な単位表現の除去処理フローの詳細の一例を示す図である。

【図5】不要な単位表現の除去処理フローの詳細の一例を示す図である。

【図6】不要な単位表現の除去処理フローの詳細の一例を示す図である。

【図7】不要な単位表現の除去処理フローの詳細の一例を示す図である。

【図8】二項分布を示す図である。

【図9】不要な単位表現の除去処理フローの詳細の一例を示す図である。

【図10】不要な単位表現の除去処理フローの詳細の一例を示す図である。

#### 【符号の説明】

- 1 統計的検定を利用した質問応答システム
- 11 不要共起語除去システム
- 12 質問文入力部
- 13 キーワード抽出/情報検索部
- 14 解答抽出部
- 15 解答出力部
- 16 データベース
- 17 表現推測部
- 18 単位表現データベース

10

20

30

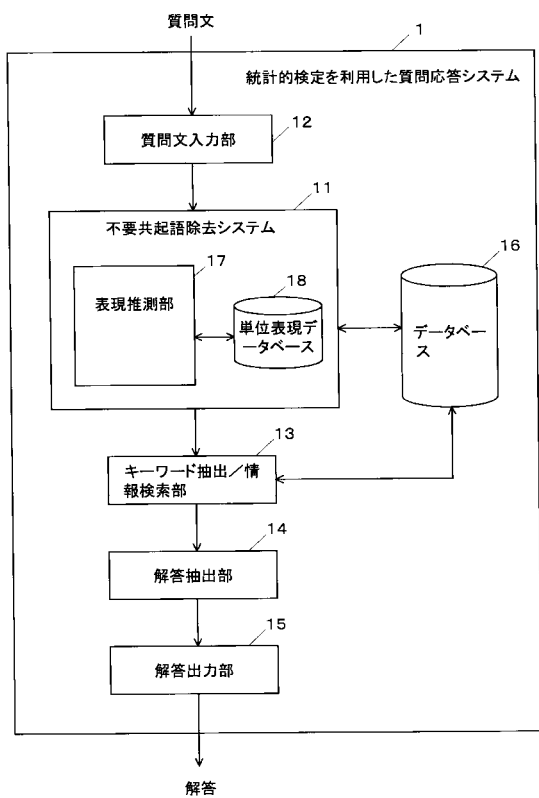
40

50

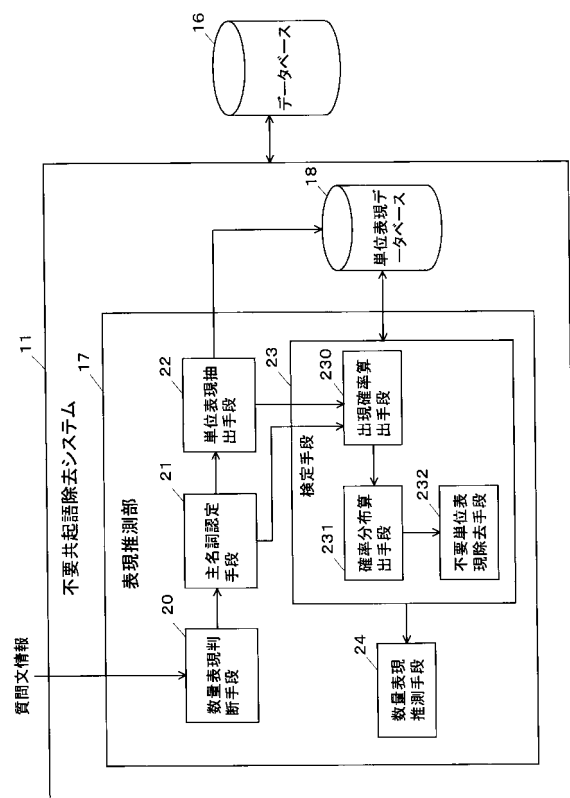


- 2 0 数量表現判断手段
- 2 1 主名詞認定手段
- 2 2 単位表現抽出手段
- 2 3 検定手段
- 2 4 数量表現推測手段
- 2 3 0 出現確率算出手段
- 2 3 1 確率分布算出手段
- 2 3 2 不要単位表現除去手段

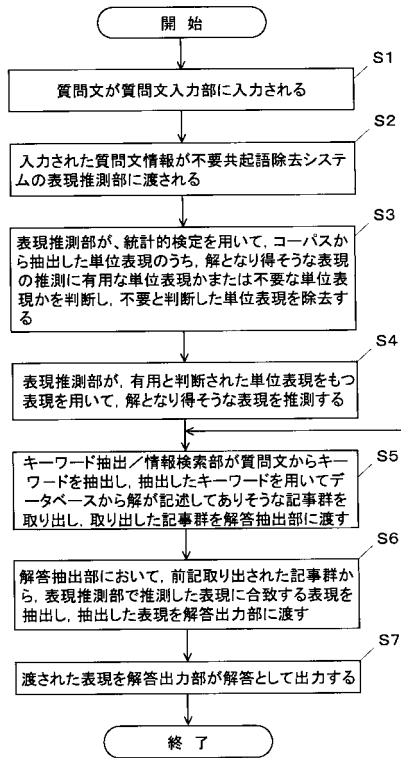
【 図 1 】



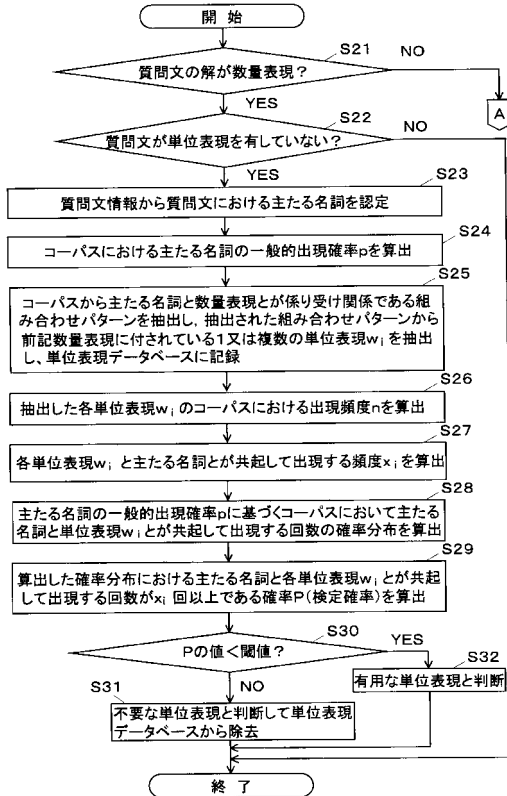
【 図 2 】



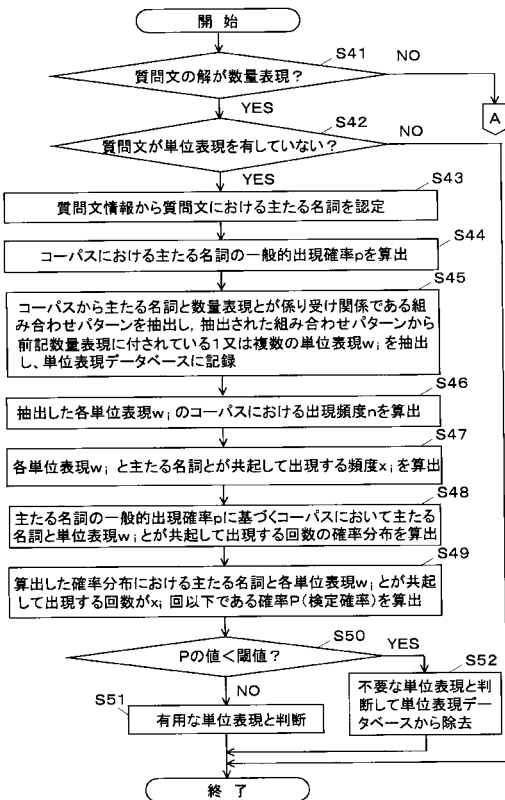
【 図 3 】



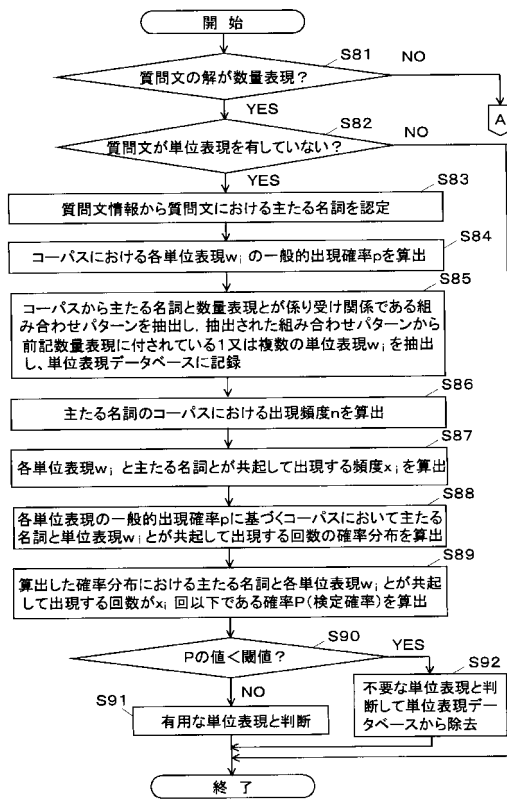
【 図 4 】



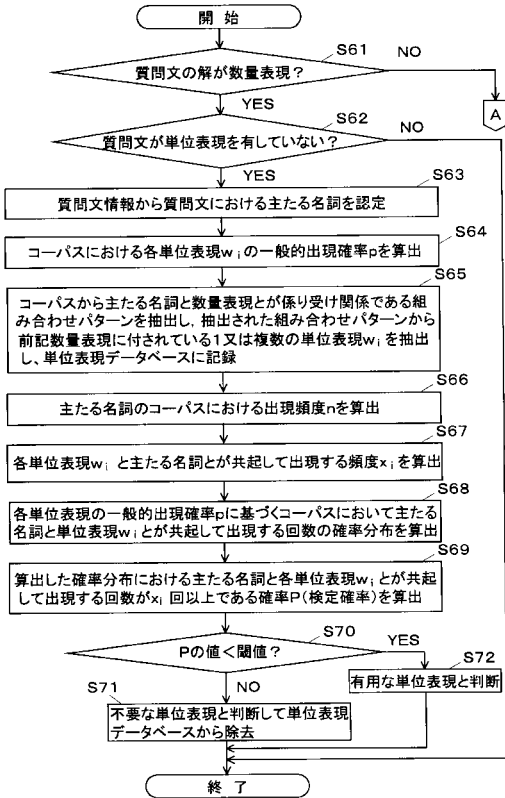
【 図 5 】



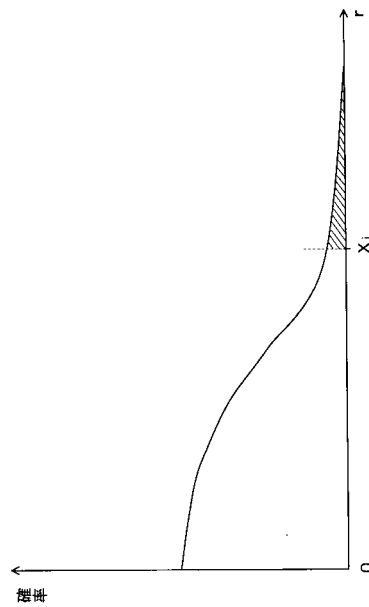
【 図 6 】



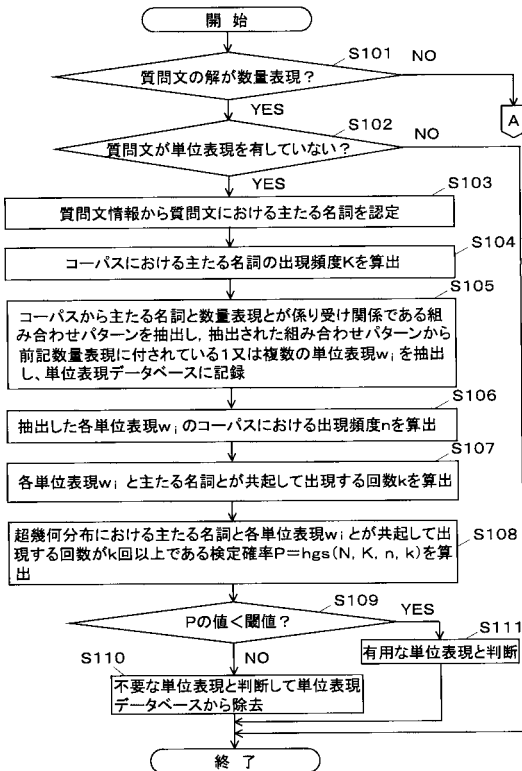
【 図 7 】



【 図 8 】



【 図 9 】



【 図 10 】

