

(19) 日本国特許庁(JP)

(12) 特 許 公 報(B2)

(11) 特許番号

特許第6165657号  
(P6165657)

(45) 発行日 平成29年7月19日 (2017.7.19)

(24) 登録日 平成29年6月30日 (2017.6.30)

(51) Int.Cl. F I  
**G 0 6 F 17/30 (2006.01)**  
 G 0 6 F 17/30 1 7 0 A  
 G 0 6 F 17/30 3 5 0 C  
 G 0 6 F 17/30 2 1 0 A

請求項の数 10 (全 23 頁)

(21) 出願番号	特願2014-58246 (P2014-58246)	(73) 特許権者	000003078 株式会社東芝
(22) 出願日	平成26年3月20日 (2014.3.20)		東京都港区芝浦一丁目1番1号
(65) 公開番号	特開2015-184749 (P2015-184749A)	(74) 代理人	110002147 特許業務法人酒井国際特許事務所
(43) 公開日	平成27年10月22日 (2015.10.22)	(72) 発明者	中田 康太 東京都港区芝浦一丁目1番1号 株式会社東芝内
審査請求日	平成28年8月30日 (2016.8.30)	(72) 発明者	蟻生 政秀 東京都港区芝浦一丁目1番1号 株式会社東芝内
		審査官	樋口 龍弥

最終頁に続く

(54) 【発明の名称】 情報処理装置、情報処理方法およびプログラム

(57) 【特許請求の範囲】

【請求項1】

複数の候補文書から言語モデルの学習に用いる文書を選択する情報処理装置であって、前記言語モデルが利用される目的に合致した目的文書について、それぞれのトピックに対する文書の関連の強さを表すトピック特徴量を算出する第1特徴量算出部と、

前記複数の候補文書のそれぞれについて、前記トピック特徴量を算出する第2特徴量算出部と、

前記複数の候補文書のそれぞれの前記トピック特徴量について、前記目的文書の前記トピック特徴量との類似度を算出する類似度算出部と、

前記類似度が基準値より大きい候補文書を、前記言語モデルの学習に用いる文書として

選択する選択部と、

【請求項2】

を備える情報処理装置。

トピック毎に、単語と、前記単語のトピックとの関連の強さを表すスコアとのペアの集合を含むトピック情報を取得するトピック情報取得部をさらに備え、

前記第1特徴量算出部および前記第2特徴量算出部は、前記トピック情報に基づき、前記トピック特徴量を算出する

請求項1に記載の情報処理装置。

【請求項3】

前記第1特徴量算出部および前記第2特徴量算出部は、トピック毎に、対象の文書に含

10

20

まれる単語のスコアを累積して、前記トピック特徴量を算出する

請求項 2 に記載の情報処理装置。

【請求項 4】

選択された前記候補文書に基づき、前記言語モデルを学習する学習部  
をさらに備える請求項 1 に記載の情報処理装置。

【請求項 5】

前記トピック情報取得部は、前記複数の候補文書を用いて前記トピック情報を生成する  
請求項 2 に記載の情報処理装置。

【請求項 6】

前記トピック情報取得部は、異なるトピック数の複数の前記トピック情報を生成し、生  
成した複数の前記トピック情報に基づき、前記目的文書の複数の前記トピック特徴量を算  
出し、算出した複数の前記トピック特徴量に基づき、生成した複数の前記トピック情報の  
うちの 1 つの前記トピック情報を選択する

請求項 5 に記載の情報処理装置。

【請求項 7】

前記トピック情報取得部は、品詞群毎に前記トピック情報を生成し、  
前記第 1 特徴量算出部および前記第 2 特徴量算出部は、前記品詞群毎の前記トピック情  
報に基づき、前記品詞群毎の前記トピック特徴量を算出する

請求項 5 に記載の情報処理装置。

【請求項 8】

前記目的文書と内容が異なり前記言語モデルの学習の基準となる  
学習対象の言語モデルと類似した用途で用いられる言語モデルを学習するための類似目  
的の文書に対する、品詞群毎の前記トピック特徴量を算出する第 3 特徴量算出部をさらに備  
え、

前記類似度算出部は、

前記複数の候補文書のそれぞれの第 1 の品詞群に関する前記トピック特徴量に対して、  
前記目的文書の前記第 1 の品詞群に関する前記トピック特徴量との第 1 の類似度を算出し、

前記複数の候補文書のそれぞれの第 2 の品詞群に関する前記トピック特徴量に対して、  
前記類似目的文書の前記第 2 の品詞群に関する前記トピック特徴量との第 2 の類似度を算  
出し、

前記選択部は、前記第 1 の類似度が第 1 の基準値より大きく、且つ、前記第 2 の類似度  
が第 2 の基準値より大きい候補文書を、前記言語モデルの学習に用いる文書として選択す  
る

請求項 7 に記載の情報処理装置。

【請求項 9】

複数の候補文書から言語モデルの学習に用いる文書を選択する情報処理方法であって、  
前記言語モデルが利用される目的に合致した目的文書について、それぞれのトピックに  
対する文書の関連の強さを表すトピック特徴量を算出する第 1 特徴量算出ステップと、

前記複数の候補文書のそれぞれについて、前記トピック特徴量を算出する第 2 特徴量算  
出ステップと、

前記複数の候補文書のそれぞれの前記トピック特徴量について、前記目的文書の前記ト  
ピック特徴量との類似度を算出する類似度算出ステップと、

前記類似度が基準値より大きい候補文書を、前記言語モデルの学習に用いる文書として  
選択する選択ステップと、

を実行する情報処理方法。

【請求項 10】

コンピュータを、複数の候補文書から言語モデルの学習に用いる文書を選択する情報処  
理装置として機能させるためのプログラムであって、

前記情報処理装置は、

10

20

30

40

50

前記言語モデルが利用される目的に合致した目的文書について、それぞれのトピックに対する文書の関連の強さを表すトピック特徴量を算出する第1特徴量算出部と、

前記複数の候補文書のそれぞれについて、前記トピック特徴量を算出する第2特徴量算出部と、

前記複数の候補文書のそれぞれの前記トピック特徴量について、前記目的文書の前記トピック特徴量との類似度を算出する類似度算出部と、

前記類似度が基準値より大きい候補文書を、前記言語モデルの学習に用いる文書として選択する選択部と、

を備えるプログラム。

【発明の詳細な説明】

10

【技術分野】

【0001】

本発明の実施形態は、情報処理装置、情報処理方法およびプログラムに関する。

【背景技術】

【0002】

コンピュータおよびインターネット環境の普及により、大量の文書が電子化され蓄積されている。このような電子化された大量の文書を用いて、音声認識等の技術に利用される言語モデルを学習することができる。例えばウェブ上で公開されている大量の文書を用いて、一般的な用途に利用される言語モデルを学習することにより、その言語モデルの性能を向上させることができる。しかし、ある特定の目的に利用される言語モデルをウェブ上で公開されている大量の文書を用いて学習しても、特定の目的以外に関する文書が多量に含まれるので、性能を大幅に向上させることはできない。

20

【0003】

ある特定の目的に利用される言語モデルの性能を向上させるには、特定の目的に関する文書（目的文書）のみを用いて言語モデルを学習すればよい。例えば、特定の目的がコールセンターにおける音声認識である場合、コールセンターにおけるオペレータのやり取りの音声を書き起こした文書を用いて言語モデルを学習すれば、その特定の目的に利用される言語モデルの性能を向上させることができる。

【0004】

ところで、このような方法は、十分な量の目的文書を用いて学習しなければ、多様な表現に対応した言語モデルとすることができない。しかし、特定の目的に関する文書を数多く収集することは困難である。例えば、音声を書き起こして文書化する作業は、経済的および時間的なコストが大きく、十分な量の目的文書を得ることは困難である。

30

【先行技術文献】

【特許文献】

【0005】

【特許文献1】特開2009-238235号公報

【発明の概要】

【発明が解決しようとする課題】

【0006】

40

発明が解決しようとする課題は、言語モデルを学習するために用いられる文書を、目的外の文書を大量に含む複数の候補文書から適切に選択することにある。

【課題を解決するための手段】

【0007】

実施形態の情報処理装置は、複数の候補文書から言語モデルの学習に用いる文書を選択する。前記情報処理装置は、第1特徴量算出部と、第2特徴量算出部と、類似度算出部と、選択部と、を備える。前記第1特徴量算出部は、前記言語モデルが利用される目的に合致した目的文書について、それぞれのトピックに対する文書の関連の強さを表すトピック特徴量を算出する。前記第2特徴量算出部は、前記複数の候補文書のそれぞれについて、前記トピック特徴量を算出する。前記類似度算出部は、前記複数の候補文書のそれぞれの

50

前記トピック特徴量について、前記目的文書の前記トピック特徴量との類似度を算出する。前記選択部は、前記類似度が基準値より大きい候補文書を、前記言語モデルの学習に用いる文書として選択する。

【図面の簡単な説明】

【0008】

【図1】第1実施形態に係る情報処理装置の構成を示す図。

【図2】トピック数が50個のトピック情報の一例を示す図。

【図3】第1実施形態に係る情報処理装置の処理フローを示す図。

【図4】目的文書の第1例を示す図。

【図5】候補文書の第1例を示す図。

10

【図6】候補文書の第2例を示す図。

【図7】候補文書の第3例を示す図。

【図8】トピック特徴量の算出フローを示す図。

【図9】単語の一致度の高い文書の一例を示す図。

【図10】トピック数が10個のトピック情報の一例を示す図。

【図11】トピック数が200個のトピック情報の一例を示す図。

【図12】トピック情報を選択するための処理フローを示す図。

【図13】第2変形例に係るトピック情報の一例を示す図。

【図14】第2実施形態に係る情報処理装置の構成を示す図。

【図15】第2実施形態に係る情報処理装置の処理フローを示す図。

20

【図16】目的文書の第2例を示す図。

【図17】類似目的文書の一例を示す図。

【図18】第1の品詞群のトピック情報の一例を示す図。

【図19】第2の品詞群のトピック情報の一例を示す図。

【図20】情報処理装置のハードウェア構成を示す図。

【発明を実施するための形態】

【0009】

(第1の実施形態)

図1は、第1実施形態に係る情報処理装置10の構成を示す図である。図2は、トピック数が50個のトピック情報の一例を示す図である。

30

【0010】

情報処理装置10は、ウェブ上等の複数の候補文書から言語モデルの学習に用いる文書を選択し、選択した候補文書を用いて言語モデルを学習する。情報処理装置10は、目的文書格納部21と、候補コーパス格納部22と、トピック情報取得部23と、第1特徴量算出部24と、第2特徴量算出部25と、類似度算出部26と、選択部27と、学習部28とを備える。

【0011】

目的文書格納部21は、学習対象の言語モデルが利用される目的に合致した文書(目的文書)を格納する。目的文書は、一例として、ユーザにより手動で選択される。学習対象の言語モデルがコールセンターにおける音声認識に利用される場合には、目的文書は、一例として、コールセンターにおけるオペレータの音声を書き起こしたテキストである。

40

【0012】

候補コーパス格納部22は、言語モデルの学習に用いる文書の候補となる複数の文書(候補文書)を格納する。複数の候補文書は、一例として、ウェブから収集した大量のテキストである。複数の候補文書には、例えば、ニュースサイトの記事、および、掲示板に書き込まれたコメント等の、多様な目的で用いられる文書が含まれ、言語モデルが利用される目的以外で用いられる文書も含まれる。候補コーパス格納部22は、情報処理装置10内に設けられるのではなく、ネットワーク上のサーバに設けられていてもよいし、複数のサーバに分散して設けられていてもよい。

【0013】

50

トピック情報取得部 23 は、トピック情報を取得する。トピック情報は、図 2 に示すような、トピック毎に、単語とスコアとのペアの集合を含む。

【0014】

トピックとは、文書で述べられている中心的な対象（テーマ）およびその文書の発話のスタイル等の特徴をいう。1つの文書に複数のトピックが含まれていてもよい。例えば、図 2 のトピック番号 # 1 は、デジタル家庭電化製品のトピックを表す。また、図 2 のトピック番号 # 2 は、食品に関するトピックを表す。さらに、トピック情報は、例えば、丁寧な発話スタイルを表すトピック、および、書き言葉のスタイル（書く場合に用いるスタイル）を表すトピックを含んでもよい。

【0015】

トピック情報におけるそれぞれのトピックに属する単語は、そのトピックに関連する単語であって、そのトピックに関する文書に含まれる可能性がある。また、トピック情報に含まれるそれぞれの単語は、スコアとペアとなっている。スコアは、その単語が属するトピックとの関連の強さを表す。本実施形態においては、スコアは、大きいほど、対するトピックとの関連が強いことを表す。

【0016】

なお、トピック情報は、1つの単語が、複数のトピックに属していてもよい。また、トピック情報に含まれるトピックの数は、何個であってもよい。

【0017】

トピック情報は、一例として、ユーザが複数のトピックを設定し、ユーザがそれぞれのトピックに関する単語を収集することにより、生成される。また、トピック情報は、一例として、ユーザが複数のトピックを設定し、ユーザがトピック毎に関連する文書とを準備し、コンピュータが準備した複数の文書内の単語の頻度を算出することにより、生成される。

【0018】

また、トピック情報取得部 23 は、例えば、下記の文献に記載されているような教師無しトピック分析技術により、トピック情報を自動で生成してもよい。

Blei, David M., Andrew Y. Ng, and Michael I. Jordan. "Latent dirichlet allocation." the Journal of machine Learning research 3 (2003): P. 993-1022.

【0019】

この方法では、まず、ユーザがトピック数  $N$  を設定する。そして、トピック情報取得部 23 は、大量で多様な文書を解析して、 $N$  個のトピックに分割されたトピック情報を生成する。この方法によれば、トピック情報取得部 23 は、トピックに関する事前の知識を用いずにトピック情報を生成することができる。

【0020】

第 1 特徴量算出部 24 は、トピック情報に基づいて、目的文書格納部 21 に格納された目的文書に対するトピック特徴量を算出する。トピック特徴量は、それぞれのトピックに対する、その文書の関連の強さを表す。本実施形態では、トピック特徴量は、下記の数 1 に示されるような、ベクトル（配列）により表される。

【数 1】

$$\vec{T}(t) = (T_1, T_2, \dots, T_{49}, T_{50}) = (0.74, 0.03, \dots, 0.06, 0.65)$$

【0021】

ベクトルで表されたトピック特徴量は、トピック情報に含まれるトピックの数分の成分（例えば、 $T_1, T_2, \dots, T_{49}, T_{50}$ ）を含む。トピック特徴量に含まれるそれぞれの成分は、トピック情報に含まれるそれぞれのトピックに一対一で対応する。それぞれの成分は、対応するトピックに対する、その文書の関連の強さを表す。例えば、数 1 の成分  $T_1$  は、図 2 に示すトピック情報におけるトピック番号 # 1 のトピックに対する、文書の関連の強さを表す。

10

20

30

40

50

## 【 0 0 2 2 】

このようなトピック特徴量は、その文書のトピック毎の関連の強さの分布を表している。なお、トピック特徴量のより詳細な算出方法については、後述の図 8 を参照して説明する。

## 【 0 0 2 3 】

第 2 特徴量算出部 2 5 は、トピック情報に基づいて、候補コーパス格納部 2 2 に格納されたそれぞれの候補文書に対するトピック特徴量を算出する。なお、候補文書に対するトピック特徴量は、目的文書に対するトピック特徴量と、同一の形式であり、同一の算出方法で算出される。

## 【 0 0 2 4 】

類似度算出部 2 6 は、複数の候補文書のそれぞれのトピック特徴量に対する、目的文書のトピック特徴量との類似度を算出する。すなわち、類似度算出部 2 6 は、複数の候補文書のそれぞれのトピック毎の関連の強さの分布が、目的文書のトピック毎の関連の強さの分布とどれだけ類似しているかを算出する。

## 【 0 0 2 5 】

本実施形態においては、類似度算出部 2 6 は、ベクトルで表されたトピック特徴量の内積を演算することにより、類似度を算出する。すなわち、類似度算出部 2 6 は、候補文書に対するトピック特徴量に含まれるそれぞれの成分と、目的文書に対するトピック特徴量の対応する成分とを乗算し、乗算結果を全て加算した値を類似度として算出する。

## 【 0 0 2 6 】

選択部 2 7 は、複数の候補文書のうち、類似度が基準値より大きい候補文書を、言語モデルの学習に用いる文書として選択する。ここで、基準値は、ユーザが設定した値であってもよい。また、基準値は、複数の候補文書の類似度に基づき算出された値であってもよい。例えば、基準値は、複数の候補文書の類似度の平均値、または、複数の候補文書の類似度の最大値から一定量小さい値等であってもよい。

## 【 0 0 2 7 】

学習部 2 8 は、選択部 2 7 により選択された候補文書に基づき、言語モデルを学習する。学習部 2 8 は、一例として、n グラム言語モデルを一般的な公知技術を用いて学習する。

## 【 0 0 2 8 】

図 3 は、第 1 実施形態に係る情報処理装置 1 0 の処理フローを示す図である。なお、本フローでは、家庭電化製品のメーカーのコールセンターの音声認識に用いられる言語モデルを学習する例を説明する。また、本フローでは、図 2 で示したトピック情報を用いる例を説明する。

## 【 0 0 2 9 】

処理に先立って、予めユーザにより目的文書が目的文書格納部 2 1 に格納される。目的文書格納部 2 1 は、一例として、図 4 に示されるような、テレビジョン受像機（テレビとも称する。）のリモートコントローラ（リモコンとも称する。）についての問い合わせに対する応答音声を書き起こしたテキストを、目的文書として格納する。

## 【 0 0 3 0 】

また、処理に先立って、情報処理装置 1 0 は、ウェブ等から複数の候補文書を取得し、候補コーパス格納部 2 2 に格納する。候補コーパス格納部 2 2 は、一例として、図 5、図 6 および図 7 に示されるような、候補文書を格納する。なお、図 5 に示される候補文書 C<sub>{n1}</sub> は、家庭電化製品のメーカーのコールセンターに対する、DVD レコーダについての問い合わせ音声を書き起こしたテキストである。図 6 に示される候補文書 C<sub>{n2}</sub> は、テレビの調子がおかしいとのウェブ上での書き込みのテキストである。図 7 に示される候補文書 C<sub>{n3}</sub> は、食品メーカーのコールセンターに対する、アレルギー源に対する問い合わせの音声を書き起こしたテキストである。

## 【 0 0 3 1 】

まず、ステップ S 1 1 において、トピック情報取得部 2 3 は、トピック情報を生成する

10

20

30

40

50

。トピック情報取得部 2 3 は、予め保存されているトピック情報を取得してもよい。

【 0 0 3 2 】

続いて、ステップ S 1 2 において、第 1 特徴量算出部 2 4 は、トピック毎に、目的文書に含まれる単語のスコアを累積して、目的文書のトピック特徴量を算出する。具体的には、第 1 特徴量算出部 2 4 は、図 8 のステップ S 2 1 からステップ S 2 9 に示す手順で、目的文書のトピック特徴量を算出する。

【 0 0 3 3 】

図 8 のステップ S 2 1 において、第 1 特徴量算出部 2 4 は、トピック特徴量を初期化する。本例においては、下記の数 2 に示すように、トピック特徴量に含まれる全ての成分を 0 . 0 に初期化する。

【数 2】

$$\vec{T}(t) = (T_1, T_2, \dots, T_{49}, T_{50}) = (0.0, 0.0, \dots, 0.0, 0.0)$$

【 0 0 3 4 】

続いて、第 1 特徴量算出部 2 4 は、対象の文書に含まれる全ての単語毎に、ステップ S 2 3 からステップ S 2 7 までの処理を繰り返して実行する（ステップ S 2 2 とステップ S 2 8 との間のループ処理）。第 1 特徴量算出部 2 4 は、一例として、対象の文書の先頭の単語から最後の単語まで、1 つずつ単語を選択して、ステップ S 2 3 からステップ S 2 7 の処理を実行する。

【 0 0 3 5 】

単語毎のループ処理において、第 1 特徴量算出部 2 4 は、さらに、トピック情報に示されたトピック毎に、ステップ S 2 4 からステップ S 2 6 の処理を繰り返して実行する（ステップ S 2 3 とステップ S 2 7 との間のループ処理）。第 1 特徴量算出部 2 4 は、一例として、トピック情報のトピック番号 # 1 からトピック番号 # 5 0 まで順次にトピックを選択して、ステップ S 2 4 からステップ S 2 6 の処理を実行する。

【 0 0 3 6 】

トピック毎のループ処理において、まず、ステップ S 2 4 において、第 1 特徴量算出部 2 4 は、選択した単語が、トピック情報における対象のトピックの単語の集合に含まれるか否かを判断する。含まれない場合には（ステップ S 2 4 の No）、第 1 特徴量算出部 2 4 は、処理をステップ S 2 7 に進める。含まれる場合には（ステップ S 2 4 の Yes）、第 1 特徴量算出部 2 4 は、処理をステップ S 2 5 に遷移させる。

【 0 0 3 7 】

ステップ S 2 5 において、第 1 特徴量算出部 2 4 は、トピック情報における対象のトピックの単語の集合から、選択した単語に対応する（ペアとなる）スコアを取得する。続いて、ステップ S 2 6 において、第 1 特徴量算出部 2 4 は、取得したスコアにより、トピック特徴量の対応する成分を更新する。第 1 特徴量算出部 2 4 は、一例として、トピック特徴量の対応する成分に、取得したスコアを加算する。

【 0 0 3 8 】

例えば、ループ処理の対象の単語が「テレビ」であり、ループ処理の対象のトピックがトピック番号 # 1 であるとする。この場合、トピック番号 # 1 の単語の集合の中の「テレビ」が存在する。従って、第 1 特徴量算出部 2 4 は、トピック番号 # 1 の「テレビ」に対応するスコア（0 . 1 1）を、トピック特徴量の 1 番目の成分  $T_1$  に加算する。下記の数 3 は、初期化したトピック特徴量に対して、「テレビ」に対応するスコア（0 . 1 1）を加算した後のトピック特徴量である。

【数 3】

$$\vec{T}(t) = (T_1, T_2, \dots, T_{49}, T_{50}) = (0.11, 0.0, \dots, 0.0, 0.0)$$

【 0 0 3 9 】

10

20

30

40

50

第1特徴量算出部24は、ステップS26の処理が終了すると、処理をステップS27に進める。ステップS27において、全てのトピックについて、まだ、ステップS24からステップS26の処理を終了していない場合には、第1特徴量算出部24は、処理をステップS23に戻して次のトピックについて処理を繰り返す。終了した場合には、第1特徴量算出部24は、処理をステップS28に進める。

【0040】

ステップS28において、全ての単語について、まだ、ステップS23からステップS27の処理を終了していない場合には、第1特徴量算出部24は、処理をステップS22に戻して、次の単語について処理を繰り返す。終了した場合には、第1特徴量算出部24は、処理をステップS29に進める。

10

【0041】

下記の数4は、全ての単語についての更新処理が完了した後のトピック特徴量である。本例では、目的文書にトピック番号#1に属する単語が多く含まれているので、 $T_1$ の値が他の成分より大きくなっている。

【数4】

$$\vec{T}(t) = (T_1, T_2, \dots, T_{49}, T_{50}) = (2.5, 0.1, \dots, 0.2, 2.2)$$

【0042】

ステップS29において、第1特徴量算出部24は、トピック特徴量を正規化する。本例では、下記の数5に示される演算によりトピック特徴量を正規化する。すなわち、第1特徴量算出部24は、それぞれの成分 $T_i$ を、全ての成分の二乗平均で除算することにより、トピック特徴量を正規化する。

20

【数5】

$$T_i = \frac{T_i}{\sqrt{\sum_{i=1}^{50} T_i^2}}$$

【0043】

下記の数6は、目的文書に対する正規化後のトピック特徴量を示す。

30

【0044】

【数6】

$$\vec{T}(t) = (T_1, T_2, \dots, T_{49}, T_{50}) = (0.74, 0.03, \dots, 0.06, 0.65)$$

本例において、正規化後のトピック特徴量は、各成分の二乗和が1となる。このように正規化することにより、トピック特徴量は、対象の文書が何れのトピックと関連性が強いのかを表すことができる。なお、数6のトピック特徴量は、成分 $T_3 \sim T_{48}$ までは0.0である。従って、本実施形態において、目的文書は、トピック番号#1とトピック番号#50のトピックとの関連性が強い。

40

【0045】

第1特徴量算出部24は、以上のように目的文書に対するトピック特徴量を算出する。

【0046】

図3に戻る。続いて、情報処理装置10は、候補コーパス格納部22に格納されている候補文書毎に、ステップS14からステップS17の処理を繰り返して実行する(ステップS13とステップS18との間のループ処理)。

【0047】

候補文書毎のループ処理では、まず、ステップS14において、第2特徴量算出部25は、トピック毎に、対象の文書に含まれる単語のスコアを累積して、候補文書のトピック特徴量を算出する。具体的には、第2特徴量算出部25は、図8のステップS21からス

50

ステップ S 2 9 に示した手順で、候補文書のトピック特徴量を算出する。

【 0 0 4 8 】

下記の数 7 は、候補文書  $C_{n1}$ 、候補文書  $C_{n2}$  および候補文書  $C_{n3}$  に対するトピック特徴量を示す。

【数 7】

$$\vec{T}(c_{n1}) = (0.70, 0.01, \dots, 0.04, 0.70)$$

$$\vec{T}(c_{n2}) = (0.71, 0.02, \dots, 0.69, 0.02)$$

$$\vec{T}(c_{n3}) = (0.01, 0.68, \dots, 0.09, 0.68)$$

10

【 0 0 4 9 】

なお、数 7 に示すトピック特徴量は、成分  $T_3 \sim T_{48}$  までは 0.0 である。候補文書  $C_{n1}$  は、トピック番号 # 1 およびトピック番号 # 5 0 のトピックとの関連性が強い。候補文書  $C_{n2}$  は、トピック番号 # 1 およびトピック番号 # 4 9 のトピックとの関連性が強い。候補文書  $C_{n3}$  は、トピック番号 # 2 およびトピック番号 # 5 0 のトピックとの関連性が強い。

【 0 0 5 0 】

続いて、ステップ S 1 5 において、類似度算出部 2 6 は、目的文書のトピック特徴量と候補文書のトピック特徴量との類似度を算出する。本実施形態においては、類似度算出部 2 6 は、下記の数 8 に示されるように、目的文書のトピック特徴量と、候補文書のトピック特徴量との内積を演算する。

20

【数 8】

$$sim(t, c_j) = \vec{T}(t) \cdot \vec{T}(c_j)$$

【 0 0 5 1 】

下記の数 9 は、候補文書  $C_{n1}$ 、候補文書  $C_{n2}$  および候補文書  $C_{n3}$  に対する類似度を示す。

30

【数 9】

$$sim(t, c_{n1}) = 0.74 * 0.70 + 0.03 * 0.01 + 0.06 * 0.04 + 0.65 * 0.70 = 0.98$$

$$sim(t, c_{n2}) = 0.74 * 0.71 + 0.03 * 0.02 + 0.06 * 0.69 + 0.65 * 0.02 = 0.58$$

$$sim(t, c_{n3}) = 0.74 * 0.01 + 0.03 * 0.68 + 0.06 * 0.09 + 0.65 * 0.68 = 0.48$$

【 0 0 5 2 】

候補文書  $C_{n1}$  の類似度は、0.98 となる。候補文書  $C_{n2}$  の類似度は、0.58 となる。候補文書  $C_{n3}$  の類似度は、0.48 となる。目的文書および候補文書  $C_{n1}$  は、共に、トピック番号 # 1 およびトピック番号 # 5 0 のトピックとの関連性が強いので、類似度が他よりも高くなっている。

40

【 0 0 5 3 】

続いて、ステップ S 1 6 において、選択部 2 7 は、類似度が基準値より大きいかなを判断する。類似度が基準値以下である場合には (ステップ S 1 6 の No)、選択部 2 7 は、処理をステップ S 1 8 に進める。類似度が基準値より大きい場合には (ステップ S 1 6 の Yes)、選択部 2 7 は、処理をステップ S 1 7 に進める。

【 0 0 5 4 】

ステップ S 1 7 において、選択部 2 7 は、対象の候補文書を、言語モデルの学習に用いる文書として選択する。本例においては、選択部 2 7 は、基準値が 0.70 に設定されており、類似度が 0.70 より大きい候補文書  $C_{n1}$  を選択する。そして、選択部 2

50

7は、処理をステップS18に進める。

【0055】

ステップS18において、全ての候補文書について、まだステップS14からステップS17の処理を終了していない場合には、選択部27は、処理をステップS13に戻して、次の候補文書について処理を繰り返す。終了した場合には、選択部27は、処理をステップS19に進める。

【0056】

ステップS19において、学習部28は、選択された候補文書を用いて、言語モデルを学習する。そして、ステップS19の処理を終えると、情報処理装置10は、本フローを終了する。

10

【0057】

以上のように、本実施形態に係る情報処理装置10によれば、目的外の文書を大量に含む複数の候補文書から、言語モデルを学習するために適切な文書を効率良く選択することができる。特に、情報処理装置10によれば、目的文書に含まれる単語と一致する単語が比較的少ない候補文書であっても、トピックの分布が類似していれば、言語モデルを学習するために用いる文書として選択することができる。

【0058】

例えば、図4に示す目的文書と、図5に示す候補文書C<sub>{n1}</sub>とを比較すると、含まれる単語の多くが異なっており、単語毎の一致度は低い。しかし、例えば、図4に示す目的文書の「テレビ」と図5に示す候補文書C<sub>{n1}</sub>の「DVD」とは、両者ともデ

20

ジタル家庭電化製品に関連する単語として認識されるので、人間の感覚では類似すると判断される。情報処理装置10は、このような候補文書C<sub>{n1}</sub>を選択する。

【0059】

また、単語の一致度の高い文書は、ほとんどが同一の単語を用いたテキストで構成される可能性がある。例えば、図9は、図4に示す目的文書と、単語の一致度の高い候補文書の一例を示す図である。図9の候補文書は、目的文書とほぼ同様の表現で構成された文書となっている。従って、図9に示すような候補文書を用いて言語モデルを学習したとしても、多様な表現に対して脆弱な言語モデルとなってしまう。

【0060】

情報処理装置10は、目的文書および候補文書のトピック特徴量を比較して類似度を判断する。従って、情報処理装置10は、目的文書と単語の一致度が低くても、同一のトピックに属する単語が含まれる候補文書を選択することができる。例えば、図5に示す候補文書C<sub>{n1}</sub>は、図4に示す目的文書と同様に、トピック番号#1およびトピック番号#50のトピックの成分が大きいので、言語モデルを学習するための文書として選択される。従って、情報処理装置10では、人間の感覚では目的文書と類似すると判断される候補文書を適切に選択することができる。これにより、情報処理装置10によれば、目的に関する多様な表現を含む文書により言語モデルを学習することができるので、多様な表現に対して頑健な言語モデルを生成することができる。

30

【0061】

(第1変形例)

つぎに、第1実施形態の第1変形例に係る情報処理装置10について説明する。

40

【0062】

図10は、トピック数が10個のトピック情報の一例を示す図である。図11は、トピック数が200個のトピック情報の一例を示す図である。

【0063】

トピック数が少ない場合、1つのトピックには、広い範囲に関連する単語が含まれる。例えば、図10に示されるように、トピック数が10個のトピック情報には、トピック番号#1のトピックに「テレビ」「DVD」等のデジタル家庭電化製品に関連する単語に加えて、「番組」「年末」等のテレビジョン番組に関連する単語が含まれてしまう。

【0064】

50

トピック数が多い場合、1つのトピックには、狭い範囲に関連する単語が含まれる。例えば、図11に示されるように、トピック数が200個のトピック情報には、トピック番号#1のトピックとトピック番号#2のトピックとに、「テレビ」と「DVD」とが分かれて属してしまう。そして、トピック番号#1には「テレビ」に関連する単語が含まれ、トピック番号#2には「DVD」に関連する単語が含まれてしまう。

【0065】

そこで、第1変形例に係るトピック情報取得部23は、複数のトピック数Nに対してトピック情報を生成し、生成されたトピック情報の中から最も適切なトピック情報を選択する。

【0066】

図12は、適切なトピック数のトピック情報を選択するための処理フローを示す図である。

【0067】

まず、ステップS31において、トピック情報取得部23は、トピック数が異なる複数のトピック情報を生成する。本例においては、トピック情報取得部23は、トピック数N=10、N=50、N=200のトピック情報を生成する。

【0068】

続いて、ステップS32において、トピック情報取得部23は、トピック数が異なる複数のトピック情報のそれぞれに基づいて、目的文書のトピック特徴量を算出する。下記の数10は、トピック数N=10、N=50、N=200の場合のトピック情報を示す。なお、数10に示すトピック特徴量は、 $T_3$ 以降の成分の値は0.0である。

【数10】

$$\vec{T}_{10}(t) = (T_1, T_2, \dots) = (0.80, 0.04, \dots)$$

$$\vec{T}_{50}(t) = (T_1, T_2, \dots) = (0.74, 0.03, \dots)$$

$$\vec{T}_{200}(t) = (T_1, T_2, \dots) = (0.54, 0.50, \dots)$$

【0069】

トピック数N=10およびトピック数N=50のトピック情報は、「テレビ」および「リモコン」がトピック番号#1のトピックに属する。従って、トピック数N=10およびトピック数N=50のトピック情報に基づく、トピック特徴量は、トピック番号#1の成分 $T_1$ の値が大きい。

【0070】

トピック数N=200のトピック情報は、「テレビ」がトピック番号#1のトピックに属し、「リモコン」がトピック番号#2のトピックに属する。従って、トピック数N=200のトピック情報に基づく、トピック特徴量は、トピック番号#1の成分 $T_1$ とトピック番号#2の成分 $T_2$ がほぼ同等となっている。

【0071】

続いて、ステップS33において、トピック情報取得部23は、生成した複数のトピック情報のうち、含まれる最大の成分の値が、閾値以上であるトピック情報を抽出する。本例の場合、トピック数N=10のトピック情報に基づくトピック特徴量の最大の成分の値は、0.80である。また、トピック数N=50のトピック情報に基づくトピック特徴量の最大の成分の値は、0.74である。また、トピック数N=200のトピック情報に基づくトピック特徴量の最大の成分の値は、0.54である。そして、閾値を0.7とした場合、トピック情報取得部23は、閾値以上であるトピック情報として、トピック数N=10のトピック情報、および、トピック数N=50のトピック情報を抽出する。

【0072】

続いて、ステップS34において、トピック情報取得部23は、抽出したトピック情報

10

20

30

40

50

のうち、トピック数が最大となるトピック情報を選択する。本例の場合、トピック情報取得部 23 は、トピック数  $N = 50$  のトピック情報を選択する。

【0073】

第1変形例に係る情報処理装置10は、このように適切な数のトピック数に設定されたトピック情報を用いて、言語モデルを学習するための候補文書を選択する。これにより、本変形例に係る情報処理装置10によれば、より性能の良い言語モデルを学習することができる。

【0074】

(第2変形例)

つぎに、第1実施形態の第2変形例に係る情報処理装置10について説明する。図13は、第2変形例に係るトピック情報の一例を示す図である。

10

【0075】

第2変形例に係るトピック情報は、文章および発話のスタイルを表すトピックの単語の集合を含む。例えば、図13に示すトピック情報におけるトピック番号#49のトピックは、親しい友人との会話で使用されるような通常の発話スタイルで用いられる単語の集合を含む。また、図13に示すトピック情報におけるトピック番号#50のトピックは、接客等で用いられるような丁寧な発話スタイルで用いられる単語の集合を含む。

【0076】

例えば、コールセンターのオペレータは、通常、丁寧な発話スタイルの音声を発生する。従って、デジタル家庭電化製品に属する単語が含まれている文書であって、且つ、日本語において文章の語尾に用いられる「です」または「ます」等の丁寧な発話スタイルに用いられる単語を含む文書を選択することにより、コールセンターのオペレータの音声認識に用いられる言語モデルを、効率良く学習することができる。

20

【0077】

従って、第2変形例に係る情報処理装置10によれば、トピック情報が発話スタイルを表すトピックの単語の集合を含むことにより、特定の用途の言語モデルを学習するために、より適切な候補文書を選択することができる。

【0078】

(第2実施形態)

つぎに、第2実施形態に係る情報処理装置10について説明する。なお、第2実施形態に係る情報処理装置10は、第1実施形態に係る情報処理装置10と略同一の機能および構成を有する。従って、略同一の機能および構成を有する要素には同一の符号を付けて、相違点を除き詳細な説明を省略する。

30

【0079】

図14は、第2実施形態に係る情報処理装置10の構成を示す図である。第2変形例に係る情報処理装置10は、類似目的文書格納部61と、第3特徴量算出部62とをさらに備える。

【0080】

類似目的文書格納部61は、学習対象の言語モデルと類似した用途で用いられる言語モデルを学習するための文書(類似目的文書)を格納する。例えば、学習対象の言語モデルが、デジタル家庭電化製品のメーカーのコールセンターの音声認識に用いられる場合であれば、類似目的文書により学習する言語モデルは、異なる商品のメーカーのコールセンターの音声認識に用いられる。

40

【0081】

トピック情報取得部23は、含まれる単語が品詞群毎に分割されたトピック情報を取得する。トピック情報取得部23は、一例として、名詞(第1の品詞群)を含むトピック情報と、名詞以外の単語(例えば、助詞、助動詞、動詞および代名詞等の第2の品詞群)を含むトピック情報とを生成する。

【0082】

第1特徴量算出部24は、品詞群毎のトピック情報に基づき、目的文書に対する品詞群

50

毎のトピック特徴量を算出する。第1特徴量算出部24は、一例として、目的文書に対する、名詞(第1の品詞群)に関するトピック特徴量および名詞以外の単語(第2の品詞群)に関するトピック特徴量を算出する。

【0083】

第2特徴量算出部25は、品詞群毎に分割されたトピック情報に基づき、それぞれの候補文書に対する品詞群毎のトピック特徴量を算出する。第2特徴量算出部25は、一例として、候補文書に対する、名詞(第1の品詞群)に関するトピック特徴量および名詞以外の単語(第2の品詞群)に関するトピック特徴量を算出する。

【0084】

第3特徴量算出部62は、品詞群毎に分割されたトピック情報に基づき、類似目的文書に対する品詞群毎のトピック特徴量を算出する。第3特徴量算出部62は、一例として、類似目的文書に対する、名詞(第1の品詞群)に関するトピック特徴量および名詞以外の単語(第2の品詞群)に関するトピック特徴量を算出する。

10

【0085】

類似度算出部26は、第1算出部71と、第2算出部72とを有する。第1算出部71は、目的文書に対する品詞群毎のトピック特徴量、および、それぞれの候補文書に対する品詞群毎のトピック特徴量を入力する。また、第1算出部71は、第1の品詞群の指定を入力する。そして、第1算出部71は、複数の候補文書のそれぞれの第1の品詞群に関するトピック特徴量に対して、目的文書の第1の品詞群に関するトピック特徴量との第1の類似度を算出する。第1算出部71は、一例として、それぞれの候補文書の名詞(第1の品詞群)に関するトピック特徴量に対して、目的文書の名詞(第1の品詞群)に関するトピック特徴量の類似度(第1の類似度)を算出する。

20

【0086】

第2算出部72は、類似目的文書に対する品詞群毎のトピック特徴量、および、それぞれの候補文書に対する品詞群毎のトピック特徴量を入力する。また、第2算出部72は、第2の品詞群の指定を入力する。そして、第2算出部72は、複数の候補文書のそれぞれの第2の品詞群に関するトピック特徴量に対して、類似目的文書の第2の品詞群に関するトピック特徴量との第2の類似度を算出する。第2算出部72は、一例として、それぞれの候補文書の名詞以外の品詞(第2の品詞群)に関するトピック特徴量に対して、類似目的文書の名詞以外の品詞(第2の品詞群)に関するトピック特徴量の類似度(第2の類似度)を算出する。

30

【0087】

選択部27は、複数の候補文書のうち、第1の類似度が第1の基準値より大きく、且つ、第2の類似度が第2の基準値より大きい候補文書を、言語モデルの学習に用いる文書として選択する。

【0088】

ここで、第1の基準値および第2の基準値は、ユーザが設定した値であってもよい。また、第1の基準値は、複数の候補文書の第1の類似度に基づき算出された値(平均値または最大値に基づく値等)であってもよい。また、第2の基準値は、複数の候補文書の第2の類似度に基づき算出された値(平均値または最大値に基づく等)であってもよい。

40

【0089】

図15は、第2実施形態に係る情報処理装置10の処理フローを示す図である。なお、本フローでは、家庭電化製品のメーカーのコールセンターの音声認識に用いられる言語モデルを学習する例を説明する。

【0090】

処理に先立って、予めユーザにより目的文書が目的文書格納部21に格納される。目的文書格納部21は、一例として、図16に示されるような、家庭電化製品のメーカーのコールセンターのオペレータにより作成された、対話内容をまとめたレポート等のテキストを、目的文書として格納する。

【0091】

50

また、処理に先立って、情報処理装置 10 は、ウェブ等から複数の候補文書を取得して、候補コーパス格納部 22 に格納する。候補コーパス格納部 22 は、一例として、第 1 実施形態と同様の、図 5、図 6 および図 7 に示されるような、候補文書を格納する。

【0092】

また、処理に先立って、予めユーザにより類似目的文書が類似目的文書格納部 61 に格納される。類似目的文書格納部 61 は、一例として、図 17 に示されるようなテキストを類似目的文書として格納する。図 17 のテキストは、家庭電化製品とは異なる製品（食品）のメーカーのコールセンターの音声認識に用いられる言語モデルの学習に利用される文書である。

【0093】

まず、ステップ S41 において、トピック情報取得部 23 は、品詞群毎に、トピック情報を生成する。下記の数 11 は、本実施形態の品詞群の集合の一例を示す式である。

【数 11】

$$PoS = (A, B) = ([名詞], [助詞, 助動詞, 動詞, 代名詞])$$

【0094】

数 11 の式では、第 1 の品詞群 A は、名詞であることを示し、第 2 の品詞群 B は、助詞、助動詞、動詞および代名詞であることを示す。なお、トピック情報取得部 23 は、3 以上の品詞群に分割したトピック情報を生成してもよい。

【0095】

トピック情報取得部 23 は、一例として、第 1 の品詞群 A のトピック情報として、図 18 に示すようなトピック情報を生成する。また、トピック情報取得部 23 は、一例として、第 2 の品詞群 B のトピック情報として、図 19 に示すようなトピック情報を生成する。

【0096】

このように品詞群毎にトピック情報を生成することにより、例えば、名詞のトピック情報は、「デジタル家庭電化製品」（トピック番号 # A\_\_1）または「食品」（トピック番号 # A\_\_2）等のトピック毎に、名詞である単語を分類することができる。また、助詞、助動詞、動詞および代名詞のトピック情報は、「書く場合に用いるスタイル」（トピック番号 # B\_\_1）または「丁寧な発話のスタイル」（トピック番号 # B\_\_2）等の文章または発話のスタイル毎に単語を分類することができる。なお、第 1 の品詞群のトピック情報と第 2 の品詞群のトピック情報とは、トピック数が異なっていてよい。

【0097】

続いて、ステップ S42 において、第 1 特徴量算出部 24 は、品詞群毎のトピック情報に基づき、目的文書に対する品詞群毎のトピック特徴量を算出する。下記の数 12 は、目的文書に対する第 1 の品詞群 A に関するトピック特徴量、および、目的文書に対する第 2 の品詞群 B に関するトピック特徴量を示す。

【数 12】

$$\vec{T}_A(t) = (T_{A1}, T_{A2}, \dots) = (0.74, 0.03, \dots)$$

$$\vec{T}_B(t) = (T_{B1}, T_{B2}, \dots) = (0.81, 0.09, \dots)$$

【0098】

数 12 に示されるように、目的文書は、トピック番号 # A\_\_1 およびトピック番号 # B\_\_1 の値が大きいため、「デジタル家庭電化製品」および「書く場合に用いるスタイル」との関連性が高いことがわかる。

【0099】

続いて、ステップ S43 において、第 3 特徴量算出部 62 は、品詞群毎のトピック情報に基づき、類似目的文書に対する品詞群毎のトピック特徴量を算出する。下記の数 13 は、類似目的文書に対する第 1 の品詞群 A に関するトピック特徴量、および、類似目的文書

10

20

30

40

50

に対する第2の品詞群Bに関するトピック特徴量を示す。

【0100】

【数13】

$$\vec{T}_A(t') = (0.01, 0.85, \dots)$$

$$\vec{T}_B(t') = (0.10, 0.80, \dots)$$

数13に示されるように、類似目的文書は、トピック番号#A\_\_2およびトピック番号#B\_\_2の値が大きいため、「食品」および「丁寧な発話スタイル」との関連性が高いことがわかる。

10

【0101】

続いて、情報処理装置10は、候補コーパス格納部22に格納されている候補文書毎に、ステップS45からステップS49の処理を繰り返して実行する（ステップS44とステップS50との間のループ処理）。

【0102】

候補文書毎のループ処理では、まず、ステップS45において、第2特徴量算出部25は、候補文書に対する品詞群毎のトピック特徴量を算出する。下記の数14は、候補文書C\_\_{n1}、候補文書C\_\_{n2}および候補文書C\_\_{n3}に対する、第1の品詞群Aおよび第2の品詞群Bに関するトピック特徴量を示す。

20

【数14】

$$\left\{ \begin{array}{l} \vec{T}_A(c_{n1}) = (0.79, 0.01, \dots) \\ \vec{T}_B(c_{n1}) = (0.10, 0.80, \dots) \\ \vec{T}_A(c_{n2}) = (0.76, 0.06, \dots) \\ \vec{T}_B(c_{n2}) = (0.75, 0.10, \dots) \\ \vec{T}_A(c_{n3}) = (0.03, 0.84, \dots) \\ \vec{T}_B(c_{n3}) = (0.06, 0.79, \dots) \end{array} \right.$$

30

【0103】

数14に示すように、候補文書C\_\_{n1}は、トピック番号#A\_\_1およびトピック番号#B\_\_2の値が大きいため、「デジタル家庭電化製品」および「丁寧な発話スタイル」との関連性が高いことがわかる。また、候補文書C\_\_{n2}は、トピック番号#A\_\_1およびトピック番号#B\_\_1の値が大きいため、「デジタル家庭電化製品」および「書く場合に用いるスタイル」との関連性が高いことがわかる。また、候補文書C\_\_{n3}は、トピック番号#A\_\_2およびトピック番号#B\_\_2の値が大きいため、「食品」および「丁寧な発話スタイル」との関連性が高いことがわかる。

40

【0104】

続いて、ステップS46において、類似度算出部26の第1算出部71は、品詞群毎に、目的文書のトピック特徴量と候補文書のトピック特徴量との類似度（第1の類似度）を算出する。本実施形態においては、第1算出部71は、下記の数15に示されるように、第1の品詞群Aおよび第2の品詞群Bのそれぞれについて、目的文書のトピック特徴量と、候補文書のトピック特徴量との内積を演算する。

50

【数 1 5】

$$sim_A(t, c_j) = \vec{T}_A(t) \cdot \vec{T}_A(c_j)$$

$$sim_B(t, c_j) = \vec{T}_B(t) \cdot \vec{T}_B(c_j)$$

【0105】

続いて、ステップ S 4 7 において、類似度算出部 2 6 の第 2 算出部 7 2 は、品詞群毎に、類似目的文書のトピック特徴量と候補文書のトピック特徴量との類似度（第 2 の類似度）を算出する。本実施形態においては、第 1 算出部 7 1 は、下記の数 1 6 に示されるように、第 1 の品詞群 A および第 2 の品詞群 B のそれぞれについて、類似目的文書のトピック特徴量と、候補文書のトピック特徴量との内積を演算する。

10

【数 1 6】

$$sim_A(t', c_j) = \vec{T}_A(t') \cdot \vec{T}_A(c_j)$$

$$sim_B(t', c_j) = \vec{T}_B(t') \cdot \vec{T}_B(c_j)$$

【0106】

続いて、ステップ S 4 8 において、選択部 2 7 は、第 1 の類似度が第 1 の基準値（ $th_A$ ）より大きく、且つ、第 2 の類似度が第 2 の基準値（ $th_B$ ）より大きいかが否かを判断する。下記の数 1 7 は、選択部 2 7 による判断条件を示す式である。

20

【数 1 7】

$$sim_A(t, c_n) > th_A \text{ かつ } sim_B(t', c_n) > th_B$$

【0107】

条件を満たさない場合には（ステップ S 4 8 の No）、選択部 2 7 は、処理をステップ S 5 0 に進める。条件を満たす場合には（ステップ S 4 8 の Yes）、選択部 2 7 は、処理をステップ S 4 9 に進める。

【0108】

ステップ S 4 9 において、選択部 2 7 は、対象の候補文書を、言語モデルの学習に用いる文書として選択する。本例においては、選択部 2 7 は、第 1 の基準値および第 2 の基準値が 0.50 に設定されており、第 1 の類似度および第 2 の類似度が共に 0.50 より大きい候補文書  $C_{n1}$  を選択する。そして、選択部 2 7 は、処理をステップ S 5 0 に進める。

30

【0109】

ステップ S 5 0 において、全ての候補文書について、まだステップ S 4 5 からステップ S 4 9 の処理を終了していない場合には、選択部 2 7 は、処理をステップ S 4 4 に戻して、次の候補文書について処理を繰り返す。終了した場合には、選択部 2 7 は、処理をステップ S 5 1 に進める。

40

【0110】

ステップ S 5 1 において、学習部 2 8 は、選択された候補文書を用いて、言語モデルを学習する。そして、ステップ S 5 1 の処理を終えると、情報処理装置 1 0 は、本フローを終了する。

【0111】

ここで、第 2 実施形態においては、候補文書  $C_{n1}$  についての数 1 7 の条件式は、下記の通りとなる。

$$sim_A(t, C_{n1}) = 0.74 * 0.79 + 0.11 * 0.03 = 0.59、\text{ かつ、 } sim_B(t', C_{n1}) = 0.10 * 0.10 + 0.8 * 0.8 = 0.65$$

50

## 【 0 1 1 2 】

従って、候補文書  $C_{n1}$  は、第 1 の品詞群 A および第 2 の品詞群 B の両方で条件を満たすので、学習用の文書として抽出される。候補文書  $C_{n1}$  は、デジタル家庭電化製品についての丁寧な発話スタイルの文書であり、コールセンターで発話される内容と一致する。従って、情報処理装置 10 は、このような文書を用いて学習を行うことで、性能の高い言語モデルを生成することができる。

## 【 0 1 1 3 】

もし、第 1 の品詞群および第 2 の品詞群の両方に対して、目的文書との類似度を用いた場合、候補文書  $C_{n1}$  についての、第 2 の品詞群 B に関する数 17 の条件式は、 $sim_B(t, C_{n1}) = 0.15$  となる。従って、この場合、候補文書  $C_{n1}$  は、条件を満たさず、学習用の文書として選択されない。一方で、候補文書  $C_{n2}$  についての数 17 の条件式は、 $sim_A(t, C_{n2}) = 0.56$ 、 $sim_B(t, C_{n2}) = 0.65$  となる。従って、この場合、候補文書  $C_{n2}$  が学習用の文書として選択され、コールセンターで実際には発話されないような、書く場合に用いるスタイルの単語を含んだ文書が、学習用の文書として選択されてしまう。

10

## 【 0 1 1 4 】

また、もし、第 1 の品詞群および第 2 の品詞群の両方に対して、類似目的文書との類似度を用いた場合には、候補文書  $C_{n1}$  についての、第 1 の品詞群 A に関する数 17 の条件式は、 $sim_A(t', C_{n1}) = 0.11$  となる。従って、この場合、候補文書  $C_{n1}$  は、条件を満たさず、学習用の文書として選択されない。

20

## 【 0 1 1 5 】

一方で、候補文書  $C_{n3}$  についての数 17 の条件式は、 $sim_A(t', C_{n3}) = 0.71$ 、 $sim_B(t, C_{n3}) = 0.64$  となる。従って、この場合、候補文書  $C_{n3}$  が学習用の文書として選択され、異なる話題のコールセンターの発話と類似した文書が、学習用の文書として選択されてしまう。

## 【 0 1 1 6 】

このように第 2 実施形態に係る情報処理装置 10 によれば、目的文書の主要なテーマと、類似目的文書の発話スタイルが予め分かっている場合に、両文書の特徴を組み合わせ、目的に合った学習用の文書を選択することができる。

## 【 0 1 1 7 】

(ハードウェア構成)

図 20 は、実施形態に係る情報処理装置 10 のハードウェア構成の一例を示す図である。実施形態に係る情報処理装置 10 は、CPU 101 (Central Processing Unit) 等の制御装置と、ROM 102 (Read Only Memory) および RAM 103 (Random Access Memory) 等の記憶装置と、ネットワークに接続して通信を行う通信 I/F 104 と、各部を接続するバスとを備えている。

30

## 【 0 1 1 8 】

実施形態に係る情報処理装置 10 で実行されるプログラムは、ROM 102 等に予め組み込まれて提供される。また、実施形態に係る情報処理装置 10 で実行されるプログラムは、インストール可能な形式または実行可能な形式のファイルで CD-ROM (Compact Disk Read Only Memory)、フレキシブルディスク (FD)、CD-R (Compact Disk Recordable)、DVD (Digital Versatile Disk) 等のコンピュータで読み取り可能な記録媒体に記録してコンピュータプログラムプロダクトとして提供されてもよい。

40

## 【 0 1 1 9 】

さらに、実施形態に係る情報処理装置 10 で実行されるプログラムは、インターネット等のネットワークに接続されたコンピュータ上に格納され、情報処理装置 10 がネットワーク経由でダウンロードすることにより提供されてもよい。また、実施形態に係る情報処理装置 10 で実行されるプログラムは、インターネット等のネットワーク経由で提供または配布されてもよい。

## 【 0 1 2 0 】

50

実施形態に係る情報処理装置 10 で実行されるプログラムは、トピック情報取得モジュール、第 1 特徴量算出モジュール、第 2 特徴量算出モジュール、第 3 特徴量算出モジュール、類似度算出モジュール、選択モジュールおよび学習モジュールを含む構成となっており、コンピュータを上述した情報処理装置 10 の各部（トピック情報取得部 23、第 1 特徴量算出部 24、第 2 特徴量算出部 25、類似度算出部 26、第 3 特徴量算出部 62、選択部 27 および学習部 28）として機能させる。このコンピュータは、CPU 101 がコンピュータ読取可能な記憶媒体からこのプログラムを主記憶装置上に読み出して実行することができる。なお、トピック情報取得部 23、第 1 特徴量算出部 24、第 2 特徴量算出部 25、類似度算出部 26、第 3 特徴量算出部 62、選択部 27 および学習部 28 は、一部または全部がハードウェアにより構成されていてもよい。

10

## 【0121】

本発明のいくつかの実施形態を説明したが、これらの実施形態は、例として提示したものであり、発明の範囲を限定することは意図していない。これら新規な実施形態は、その他の様々な形態で実施されることが可能であり、発明の要旨を逸脱しない範囲で、種々の省略、置き換え、変更を行うことができる。これら実施形態やその変形は、発明の範囲や要旨に含まれるとともに、請求の範囲に記載された発明とその均等の範囲に含まれる。

## 【符号の説明】

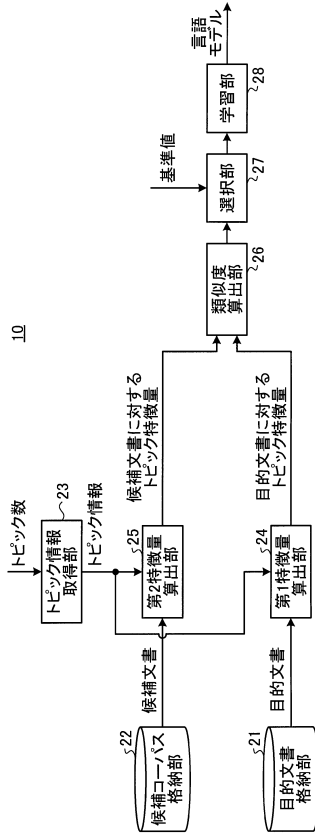
## 【0122】

- 10 情報処理装置
- 21 目的文書格納部
- 22 候補コーパス格納部
- 23 トピック情報取得部
- 24 第 1 特徴量算出部
- 25 第 2 特徴量算出部
- 26 類似度算出部
- 27 選択部
- 28 学習部
- 61 類似目的文書格納部
- 62 第 3 特徴量算出部
- 71 第 1 算出部
- 72 第 2 算出部
- 101 CPU
- 102 ROM
- 103 RAM
- 104 通信 I/F

20

30

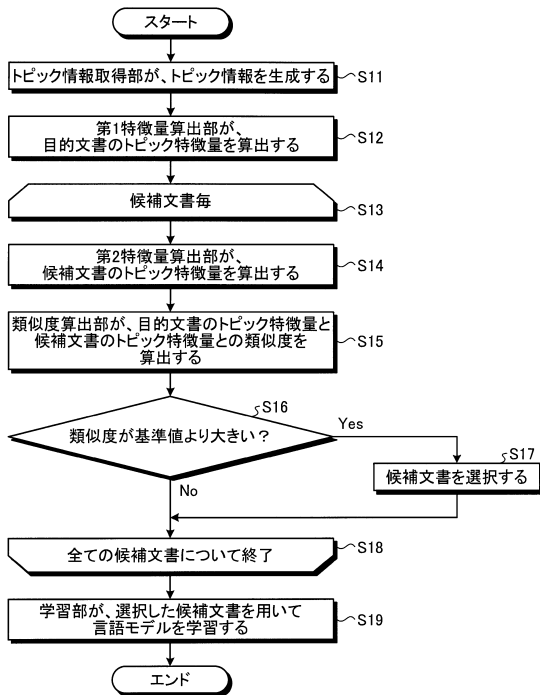
【図1】



【図2】

トピック番号#1		トピック番号#2		...		トピック番号#49		トピック番号#50	
単語	スコア	単語	スコア	単語	スコア	単語	スコア	単語	スコア
テレビ	0.11	食事	0.11	xxx	0.11	xxx	0.11	xxx	0.11
DVD	0.10	食品	0.10	xxx	0.10	xxx	0.10	xxx	0.10
リモコン	0.10	アレルギー	0.10	xxx	0.10	xxx	0.10	xxx	0.10
画面	0.09	成分	0.09	xxx	0.09	xxx	0.09	xxx	0.09
録画	0.09	ミルク	0.09	xxx	0.09	xxx	0.09	xxx	0.09
機種	0.08	添加	0.09	xxx	0.09	xxx	0.09	xxx	0.09
...	...	...	...	...	...	...	...	...	...

【図3】



【図4】

目的文書

こちらはA社コールセンターです。テレビについてのお問い合わせですか。リモコンがきかなくなったということですね。

【図5】

候補文書C<sub>[n]</sub>

DVDレコーダーについて質問があります。録画予約の仕方がわかりません。...

【図6】

候補文書C<sub>[n2]</sub>

最近、テレビの調子が悪い。  
何かリモコンもおかしいです。  
明日電気屋に相談に行こう。  
...

【図7】

候補文書C<sub>[n3]</sub>

先日購入した食品を食べたところ、  
湿疹が出てしまいました。  
何かアレルギーは入っていますか。  
...

【図9】

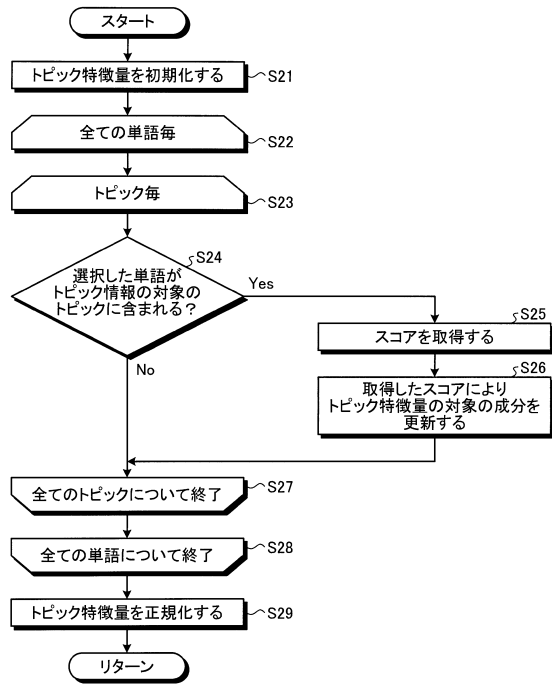
単語一致度の高い候補文書

テレビについての問い合わせです。  
リモコンがおかしいです。  
...

【図10】

トピック番号#1		...
単語	スコア	
テレビ	0.11	
番組	0.10	
冷蔵庫	0.10	
電気	0.09	
年末	0.09	
DVD	0.08	
...	...	

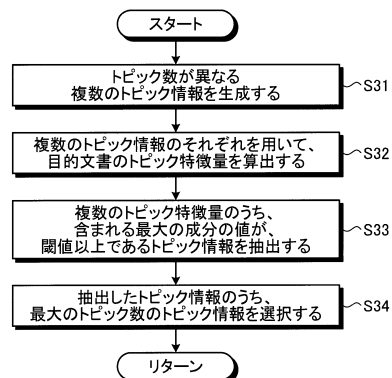
【図8】



【図11】

トピック番号#1		トピック番号#2		...
単語	スコア	単語	スコア	
テレビ	0.11	DVD	0.11	
画面	0.10	リモコン	0.10	
チャンネル	0.10	容量	0.10	
ケーブル	0.09	予約	0.09	
ベゼル	0.09	編集	0.09	
色	0.08	起動	0.09	
...	...	...	...	

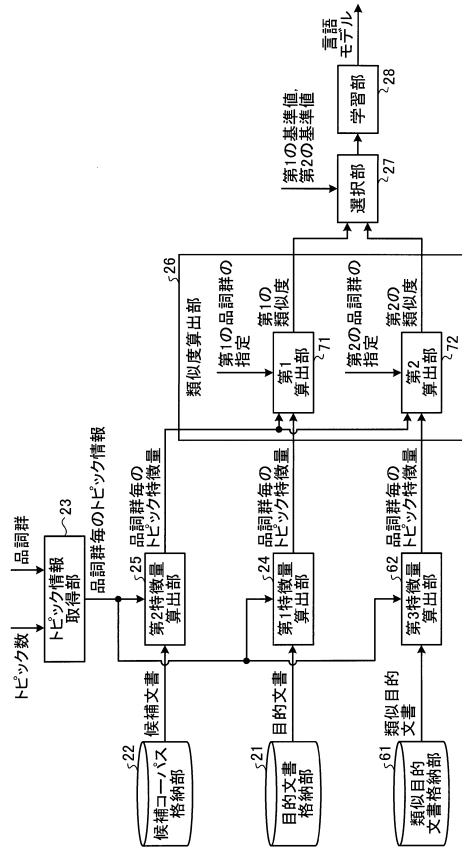
【図12】



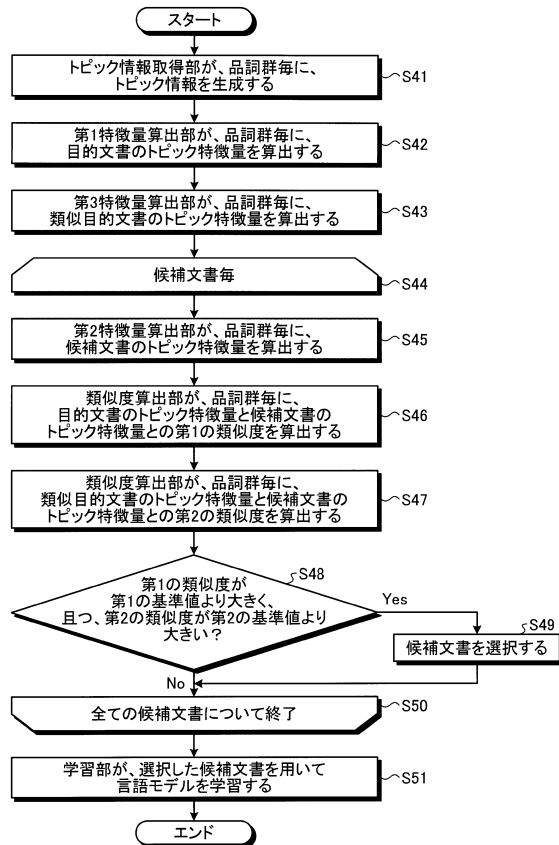
【図13】

トピック番号#1		トピック番号#2		...		トピック番号#49		トピック番号#50	
単語	スコア	単語	スコア	単語	スコア	単語	スコア	単語	スコア
テレビ	0.11	食事	0.11			って	0.11	です	0.11
DVD	0.10	食品	0.10			じゃ	0.10	ます	0.10
リモコン	0.10	アレルギー	0.10			から	0.10	私	0.10
画面	0.09	成分	0.09			けど	0.09	ましょう	0.09
録画	0.09	ミルク	0.09			それ	0.09	申し	0.09
機種	0.08	添加	0.09			みたい	0.09	先日	0.09
...	...	...	...			...	...	...	...

【図14】



【図15】



【図16】

目的文書

- ・テレビについて質問
- ・画面が見えにくい
- ・チャンネル変えても改善せず
- ...

【図17】

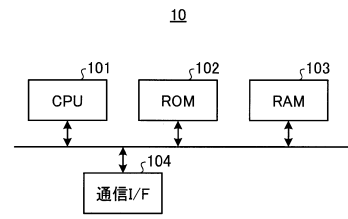
類似目的文書

B社コールセンターです。  
食品のアレルギーについての  
お問い合わせですか。  
粉ミルクの成分についてですね。  
...

【図18】

トピック番号#A.1		トピック番号#A.2		...
単語	スコア	単語	スコア	
テレビ	0.11	食事	0.11	
DVD	0.10	食品	0.10	
予約	0.10	アレルギー	0.10	
画面	0.09	成分	0.09	
録画	0.09	ミルク	0.09	
機種	0.08	添加	0.09	
...	...	...	...	

【図20】



【図19】

トピック番号#B.1		トピック番号#B.2		...
単語	スコア	単語	スコア	
ない	0.11	です	0.11	
だった	0.10	ます	0.10	
した	0.10	ましょう	0.09	
なった	0.09	申し	0.09	
する	0.09	ませ	0.09	
有り	0.09	...	...	
...	...	...	...	

---

フロントページの続き

(56)参考文献 特開2010-97318(JP,A)  
特開平4-314171(JP,A)  
特開2008-102790(JP,A)

(58)調査した分野(Int.Cl., DB名)  
G06F 17/30