

(19) 日本国特許庁(JP)

(12) 特 許 公 報(B2)

(11) 特許番号

特許第6817690号
(P6817690)

(45) 発行日 令和3年1月20日(2021.1.20)

(24) 登録日 令和3年1月4日(2021.1.4)

(51) Int.Cl.		F 1
G 0 6 F 40/279	(2020.01)	G O 6 F 40/279
G 0 6 F 16/00	(2019.01)	G O 6 F 16/00
G 0 6 F 16/30	(2019.01)	G O 6 F 16/30

請求項の数 10 (全 21 頁)

(21) 出願番号	特願2015-68461 (P2015-68461)	(73) 特許権者	000004237
(22) 出願日	平成27年3月30日 (2015. 3. 30)		日本電気株式会社
(65) 公開番号	特開2016-189089 (P2016-189089A)		東京都港区芝五丁目7番1号
(43) 公開日	平成28年11月4日 (2016. 11. 4)	(74) 代理人	100109313
審査請求日	平成30年2月15日 (2018. 2. 15)		弁理士 机 昌彦
審判番号	不服2019-15283 (P2019-15283/J1)	(74) 代理人	100124154
審判請求日	令和1年11月14日 (2019. 11. 14)		弁理士 下坂 直樹
		(72) 発明者	得地 博之
			東京都港区芝五丁目7番1号
			日本電気株式会社内
		合議体	
		審判長	渡邊 聡
		審判官	上田 智志
		審判官	中野 浩昌

最終頁に続く

(54) 【発明の名称】 抽出装置、抽出方法とそのプログラム、及び、支援装置、表示制御装置

(57) 【特許請求の範囲】

【請求項1】

複数の文を含むテキストから前記文を抽出し、前記文ごとにN個(Nは2以上の自然数)の単語をつなげたN-Gramを生成し、前記N-Gramに対し学習モデルを用いて評価する評価値を算出し、前記評価値に基づいて前記文から要約文を抽出する要約文抽出部を備える抽出装置。

【請求項2】

前記学習モデルは、複数の教師単語集合を用いて、所定の単語集合が前記教師単語集合らしいか否かを評価可能に学習されたモデルである、請求項1記載の抽出装置。

【請求項3】

前記要約文抽出部は、前記評価値に基づいて前記文ごとに教師ラベル判定寄与度を算出し、前記教師ラベル判定寄与度に応じて要約文を抽出する、請求項1又は2に記載の抽出装置。

【請求項4】

前記教師ラベル判定寄与度の算出は、前記評価値の分散値又は標準偏差値、前記評価値の最大絶対値、又は、前記評価値のノルム値のいずれかを用いる、

請求項3に記載の抽出装置。

【請求項5】

複数の文を含むテキストから前記文を抽出し、前記文ごとにN個(Nは2以上の自然数)の単語をつなげたN-Gramを生成し、前記N-Gramに対し学習モデルを用いて

評価する評価値を算出し、前記評価値に基づいて前記文から要約文を抽出する抽出方法。

【請求項 6】

複数の文を含むテキストから前記文を抽出し、前記文ごとに N 個 (N は 2 以上の自然数) の単語をつなげた N - G r a m を生成し、前記 N - G r a m に対し学習モデルを用いて評価する評価値を算出し、前記評価値に基づいて前記文から要約文を抽出することをコンピュータに実行させる抽出プログラム。

【請求項 7】

請求項 1 から 4 のいずれか 1 に記載の抽出装置と、を備え、前記抽出装置から出力された前記要約文ごとにその文中で教師単語集合らしいか否かに応じて表示を変化させる支援装置。

10

【請求項 8】

複数の文を含むテキストから抽出された前記文ごとに、学習モデルを用いて算出された、前記文から生成された N 個 (N は 2 以上の自然数) の単語をつなげた N - G r a m に対する評価値に基づいて、前記文から要約文を抽出し、前記要約文を前記評価値に基づいた順序で表示制御する表示制御部を備える表示制御装置。

【請求項 9】

複数の文を含むテキストから抽出された前記文ごとに、学習モデルを用いて算出された、前記文から生成された N 個 (N は 2 以上の自然数) の単語をつなげた N - G r a m に対する評価値に基づいて、前記文から要約文を抽出し、前記要約文を前記評価値に基づいた順序で表示制御する表示制御方法。

20

【請求項 10】

複数の文を含むテキストから抽出された前記文ごとに、学習モデルを用いて算出された、前記文から生成された N 個 (N は 2 以上の自然数) の単語をつなげた N - G r a m に対する評価値に基づいて、前記文ごとに教師ラベル判定寄与度を算出し、前記教師ラベル判定寄与度に応じて要約文を抽出し、前記要約文を前記教師ラベル判定寄与度に基づいた順序で表示制御する表示制御方法。

【発明の詳細な説明】

【技術分野】

【0001】

本発明は、抽出装置、抽出方法とそのプログラム、及び、支援装置、表示制御装置に関する、テキストからの要約文の抽出に関する。

30

【背景技術】

【0002】

昨今のビッグデータの分析需要の増加により、様々な情報分析を目的とした機械学習の教師ラベル付与の必要性が高まっている。教師ラベルは、機械学習装置に対して未知のデータを学習させる際に、そのデータがどの分類に属するか、又は、どの程度のスコアなのかを機械学習装置に教示するための情報である。ただし、教師ラベルは、装置によって自動的に判定されるのではなく、人がデータの内容を理解して教師ラベルを判定し付与する必要がある。

【0003】

40

テキストデータは、数値、又は、画像 / 映像などのデータに比べ、テキストを書いた筆者の個性、及び、意思、を表現する情報 (筆者の語彙、又は、語順 / 使用頻度の癖、および感情表現など) を豊富に含んでいるため、分析の対象として非常に有用なデータである。しかし、テキストデータは画像 / 映像データと違って一目眺めれば内容を理解できるものではなく、「読む」ことによって初めて理解することができるため、内容の理解に大きな時間を要する。また、「読む」という作業は、テキストの複雑さや長さによって作業コストが大きく上昇することから、テキストの内容理解を支援する技術が数多く発明されている。

【0004】

特許文献 1 の技術は、速読したい文書に対して文書のジャンルを特定し、ジャンルに対

50

応する決定木を選択する。一方、与えられた文書の本文中の各文について特徴を抽出する。選択された決定木と各文の特徴を照し合せ、それぞれの文について要約文か否かを決定する。要約文を強調色、非要約文を背景色で表示する。また与えられた文書の各段落の第一文目を要約文とは異なる色で表示することにより、重要箇所の抽出と表示による文書の速読支援を実現している。

【先行技術文献】

【特許文献】

【0005】

【特許文献1】特許第3652086号公報

【発明の概要】

10

【発明が解決しようとする課題】

【0006】

特許文献1では、単語の出現頻度による特徴を表すTF-IDF(Term Frequency - Inverse Document Frequency)、又は、文字数の統計的特徴を算出し、算出した単語が持つ総合的な特徴によってテキスト内の各文が要約らしいか否かを判定している。このため、単語の語順に伴う文意を反映して要約文を抽出することができない。例えば、「重要」という単語に対し、後続の単語が「である」なのか、「でない」なのかによって、文意が大きく変わる場合があり、所望の要約文の抽出ができなくなる。

【0007】

20

本発明の目的は、単語の語順に伴う文意を反映した要約文を抽出することが可能な技術を提供することにある。

【課題を解決するための手段】

【0008】

本発明の抽出装置は、複数の文を含むテキストから前記文を抽出し、前記文ごとに1以上の単語を含む単語集合を生成し、前記生成された単語集合に対し学習モデルを用いて評価する評価値を算出し、前記算出された評価値に基づいて前記文から要約文を抽出する要約文抽出部を備える。

【0009】

本発明の抽出方法は、複数の文を含むテキストから前記文を抽出し、前記文ごとに1以上の単語を含む単語集合を生成し、前記生成された単語集合に対し学習モデルを用いて評価する評価値を算出し、前記算出された評価値に基づいて前記文から要約文を抽出する。

30

【0010】

本発明の抽出プログラムは、コンピュータに、複数の文を含むテキストから前記文を抽出し、前記文ごとに1以上の単語を含む単語集合を生成し、前記生成された単語集合に対し学習モデルを用いて評価する評価値を算出し、前記算出された評価値に基づいて前記文から要約文を抽出することを実行させる。

【0011】

本発明の支援装置は、上記の抽出装置と、を備え、前記抽出装置から出力された前記要約文ごとにその文中で、前記教師単語集合らしいか否かに応じて表示を変化させる。

40

【0012】

本発明の表示制御装置は、複数の文を含むテキストから学習モデルを用いて算出された評価値に基づいて前記文から要約文を抽出し、前記要約文を前記評価値に基づいた順序で表示制御する表示制御部を備える。

【0013】

本発明の表示制御方法は、複数の文を含むテキストから学習モデルを用いて算出された評価値に基づいて前記文から要約文を抽出し、前記要約文を前記評価値に基づいた順序で表示制御する。

【発明の効果】

50

【 0 0 1 4 】

本発明の抽出装置は、単語の語順に伴う文意を反映した要約文を抽出することができる。

【 図面の簡単な説明 】

【 0 0 1 5 】

【 図 1 】 本発明の第 1 の実施形態による抽出装置の構成を示すブロック図である。

【 図 2 】 本発明の第 1 の実施形態による抽出装置の動作を示すフローチャートである。

【 図 3 】 図 2 に示す文ごとに評価値を算出するサブルーチンの動作を示すフローチャートである。

【 図 4 】 3 個の単語をつなげた単語 $N - Gram (N = 3)$ の例を説明する図である。 10

【 図 5 】 確信度の算出を説明するための図である。

【 図 6 】 図 2 に示す要約文抽出のサブルーチンの動作を示すフローチャートである。

【 図 7 】 確信度の総和によって寄与度を算出する際の問題を説明するための図である。

【 図 8 】 英語のテキストを単語 $N - Gram (N = 3)$ で処理する例を示す図である。

【 図 9 】 本発明の第 2 の実施形態による支援装置および記憶装置の構成を示すブロック図である。

【 図 1 0 】 本発明の第 2 の実施形態による支援装置の動作を示すフローチャートである。

【 図 1 1 】 図 1 0 に示す言語処理のサブルーチンの動作を示すフローチャートである。

【 図 1 2 】 図 1 0 に示す言語処理のサブルーチンの動作を示すフローチャートである。

【 図 1 3 】 第 2 の実施形態による表示装置に表示する画面表示を示す図である。 20

【 図 1 4 】 図 1 0 に示す学習のサブルーチンの動作を示すフローチャートである。

【 図 1 5 】 第 3 の実施形態による支援装置および記憶装置の構成を示すブロック図である。

【 図 1 6 】 本発明の第 4 の実施形態による表示制御装置の構成を示すブロック図である。

【 図 1 7 】 本発明の第 4 の実施形態による表示制御装置の動作を示すフローチャートである。

【 図 1 8 】 本発明の第 1 の実施形態による抽出装置、第 2、3 の実施形態による支援装置又は第 3 の実施形態による表示制御装置をコンピュータ装置で実現したハードウェア構成を示すブロック図である。

【 発明を実施するための形態 】 30

【 0 0 1 6 】

< 第 1 の実施形態 >

本発明の第 1 の実施形態である抽出装置について、図面を用いて説明する。第 1 の実施形態の抽出装置 1 0 は、テキストに教師ラベルを付与する者に対して、その教師ラベルの判定を支援する支援装置の一つの機能を提供する装置である。

【 0 0 1 7 】

図 1 は、第 1 の実施形態による抽出装置 1 0 の構成を示すブロック図である。図 1 に示すように、抽出装置 1 0 は、要約文抽出部 4 0 を備える。要約文抽出部 4 0 は、複数の文を含むテキストから文を抽出し、文ごとに 1 以上の単語を含む単語集合を生成し、生成された単語集合に対し学習モデルを用いて評価する評価値を算出し、算出された評価値に基づいて文から要約文を抽出する。 40

【 0 0 1 8 】

複数の文を含むテキストから文を抽出することの一例として、要約文抽出部 4 0 は、教師ラベルを付与するテキストである対象テキストに対し、対象テキストを構成する単語で区切った単語区切りの対象テキストを文単位に分割する。さらに、生成された単語集合に対し学習モデルを用いて評価する評価値を算出することの一例として、要約文抽出部 4 0 は、分割された文ごとに N 個の単語をつなげた単語 $N - Gram (N は 2 以上の自然数)$ を生成する。続いて要約文抽出部 4 0 は、生成された単語 $N - Gram$ に対し学習モデルを用いて教師ラベルらしさを表す確信度を算出する。さらに、算出された評価値に基づいて文から要約文を抽出することの一例として、要約文抽出部 4 0 は、算出された確信度に 50

基づいて分割された文ごとに教師ラベル判定寄与度を算出し、教師ラベル判定寄与度に応じて要約文を抽出する。教師ラベル判定寄与度については後に説明する。

【0019】

次に、本発明の第1の実施形態による抽出装置の動作について、図面を用いて説明する。図2は、第1の実施形態による抽出装置10の動作を示すフローチャートである。図2に示すように、抽出装置10は、複数の文を含むテキストから文を抽出する(S101)。具体的には、抽出装置10は、教師ラベルを付与するテキストである対象テキストに対し、対象テキストを構成する単語で区切った単語区切りの対象テキストを文単位に分割する。

【0020】

ここで、単語区切りのテキストとは、分かち書きで表現されたテキストを意味する。例えば、日本語のテキストが「お店は自宅から遠いですが、料理はとても美味しい。」である場合は、「お店は自宅から遠いですが、料理はとても美味しい。」のように単語ごとに区切られたテキストとなる。また、文単位に分割するとは、テキスト内に付された約物(句読点等)の存在及びその位置に応じて文を分けることである。例えば、前例の「お店は自宅から遠いですが、料理はとても美味しい。」という文は、読点の位置によって「お店は自宅から遠いですが、」と、「料理はとても美味しい。」という2つの文に分かれる。なお、文単位へ分割は、約物の位置以外に、次に示す単位で処理することもできる。

- ・「単語N-Gramよりも単語区切りが多い単語M-Gram($M > N$: M, Nは2以上の自然数)の単位」
- ・「K文字ごと(Kは1以上の自然数)」
- ・「行単位(改行文字)」
- ・「ページ単位(改ページコード)」
- ・「節、段落単位」

次に、抽出装置10は、文ごとに1以上の単語を含む単語集合を生成し、生成された単語集合に対し学習モデルを用いて評価する評価値を算出する(S102)。具体的には、教師ラベルらしさを表す確信度を算出し、算出された確信度に基づいて分割された文ごとに教師ラベル判定寄与度(以下、寄与度と示す。)を算出する。なお、教師ラベル判定寄与度とは、教師ラベルの付与の際に、付与する者の判定に寄与できる程度を示す値である。

【0021】

図3は、文ごとに評価値を算出するサブルーチンの動作を示すフローチャートである。図3に示すように、抽出装置10は、文ごとに単語集合を生成する(S1021)。具体的には、抽出装置10は、分割された文ごとにN個の単語をつなげた単語N-Gram(Nは2以上の自然数)を生成する。

【0022】

図4は、3個の単語をつなげた単語N-Gram($N = 3$)の例を説明する図である。図4に示すように、「私が先週予約したお店は大変好評でした。」という単語区切りされた1つの文を、単語ごとに3個の連続する単語を含む文字列に変換したものである。図4の例では、10個の単語N-Gram($N = 3$)が生成されている。

【0023】

ここで、評価値の一例である確信度とは、生成された各単語N-Gramに対して算出される教師ラベルのスコアである。よって、単語N-Gram($N = 3$)の教師ラベルのスコアとは3個の連続する単語を含む文字列が、P(ポジティブ)なのか、N(ネガティブ)なのか、その程度を表すスコアである。

当該学習モデルは、以下のように構築される。まず、学習用の教師データとして、P/N情報(ポジティブ/ネガティブ情報)が既知であるテキストが用いられる。続いて、学習モデルを生成する学習部(図示せず)は、教師データとなるテキストの単語N-Gram

10

20

30

40

50

を作成した後、単語 N-Gram ごとに単語に紐づく特徴ベクトルに置換し、学習モデルに特徴ベクトルと P/N 情報 (スコア) を教え込む。これにより、学習モデルがテキストから P/N 情報 (スコア) を判断する能力を得る。学習モデルは、例えば、サポートベクタマシン、ニューラルネットワーク、又は、ベイズ分類器のように、任意の教師あり機械学習分類器を用いて生成することができる。なお、第 1 の実施形態において、確信度を算出するための学習モデルは、確信度の算出前に予め準備されているものとする。学習モデルは、複数の教師単語集合を用いて、所定の単語集合が教師単語集合らしいか否かを評価可能に学習されたモデルであるとも言える。

【 0 0 2 4 】

次に、抽出装置 10 は、生成された単語集合に対し学習モデルを用いて評価する評価値を算出する (S 1 0 2 2)。

【 0 0 2 5 】

図 5 は、ニューラルネットワークを用いて生成した学習モデルと、生成された各単語 N-Gram とを用いた抽出装置 10 による確信度の算出を説明するための図である。抽出装置 10 は、P/N 情報が不明なテストデータ (教師ラベルを付与する対象テキスト) として、生成された各単語 N-Gram に対し学習モデルを用いて確信度を算出する。教師ラベルを付与する対象テキストとして図 4 に示す例を用いる。

【 0 0 2 6 】

図 5 に示すように、抽出装置 10 は、生成された単語 N-Gram (N = 3) ごとに単語に紐づく特徴ベクトルに置換する。次に、抽出装置 10 は、各単語 N-Gram (N = 3) ごとに置換された単語に紐づいた特徴ベクトルを、ニューラルネットワークを用いて生成した学習モデルに入力する。続いて、抽出装置 10 は、単語 N-Gram (N = 3) ごとの P/N (ポジティブ/ネガティブ) 情報のスコアを推定する。なお、確信度のスコアの範囲は、-1 から 1 まで (0.1 単位) とする。P (ポジティブ)、N (ネガティブ) の双方で現れそうな N-Gram は、「0」付近、ポジティブな文章に現れそうな N-Gram は「1」付近、ネガティブな文章に現れそうな N-Gram は「-1」付近となるように設定されている。図 5 の例では、10 個の単語 N-Gram (N = 3) ごとに、確信度 (教師ラベルのスコア) が算出される。

【 0 0 2 7 】

抽出装置 10 は、算出された評価値に基づいて文から要約文を抽出する (S 1 0 3)。

図 6 は、要約文を抽出するサブルーチンの動作を示すフローチャートである。図 6 に示すように、抽出装置 10 は、算出された評価値に基づいて文ごとに教師ラベル判定寄与度を算出する (S 1 0 3 1)。

具体的には、抽出装置 10 は、単語 N-Gram ごとに算出された確信度に基づいて、分割された文ごとに寄与度を算出する。

抽出装置 10 による寄与度の算出の一例として、次に示すバリエーションが考えられる。

- ・各単語 N-Gram における算出された確信度の分散値又は標準偏差値
- ・各単語 N-Gram における算出された確信度の最大絶対値
- ・各単語 N-Gram における算出された確信度のノルム値
- ・単語 N-Gram における算出された確信度の平均値

ここで、算出された各確信度の総和によって生じる問題について説明する。図 7 は、確信度の総和によって算出される値の一例を示す図である。図 7 に示すように、上段は、算出された 8 個の単語 N-Gram (N = 3) ごとの確信度における、ポジティブ/ネガティブ (P/N) を表し、下段は、そのスコアを表す。図 7 に示す確信度に基づき、確信度を総和だけを用いて文ごとの寄与度を算出すると、総和の合計値は、0.00 となる。すなわち、図 6 に示すように確信度としてポジティブ/ネガティブの値が極端に大きな数値であるにも関わらず、総和により、文ごとの寄与度が 0.00 となるため、後段の要約文の抽出において、その文が、重要な要約文として抽出できなくなる可能性がある。

【 0 0 2 8 】

この問題に対し、第 1 の実施形態では、寄与度の算出に、各単語 N-Gram における算出された確信度の分散値又は標準偏差を用いる。これにより、図 7 に示すように、分散

10

20

30

40

50

値が0.9、標準偏差値が0.95となり、重要な要約文を抽出することが可能となる。

【0029】

なお、「確信度の最大絶対値」を寄与度として採用することで、確信度が高い(学習モデルが自信を持って推定した)単語N-Gramが1つ以上含まれている要約文抽出も考えられる。また、これらのバリエーションの組合せによって要約文を抽出することも可能である。

【0030】

最後に、抽出装置10は、文ごとに算出された寄与度から要約文を抽出する(S1032)。要約文の抽出条件は、算出された寄与度が、所定の閾値以上である文、あるいは、算出された寄与度を降順に整列したうちの上位数十パーセントとなる文を抽出する。上記の要約文の抽出条件は一例であり、他の抽出条件でも適用可能である。

10

【0031】

また、第1の実施形態は、教師ラベルを付与するテキストとして、日本語の例を示したが、これに限られるものではなく、英語の対象テキストでも適用可能である。図8は、英語の対象テキストを単語N-Gram(N=3)で処理する例を示す図である。英語など通常、分かち書きとなっている対象テキストの場合、対象テキストを単語単位に区切る処理は不要となる。図8に示すように、抽出装置10により、文ごとに生成された単語N-Gram(N=3)ごとの確信度を算出し、算出された確信度に基づき、文ごとに教師ラベル判定寄与度を算出する。これにより、英語のテキストでも、単語の語順に伴う文意を反映した要約文を抽出することができる。

20

【0032】

第1の実施形態の抽出装置によれば、単語の語順に伴う文意を反映した要約文を抽出することが可能になる。例えば、「お店はきれいで雰囲気は悪くない。」というテキストと、「雰囲気は悪くお店はきれいでない。」というテキストでは、テキストを構成する単語は、双方とも同じになる。このため、特許文献1のように単語単位で抽出し、単語の出現頻度を用いる例では、単語の組合せで文意が変わる場合に、順序による文意を考慮することができず、所望の要約文を抽出することができない。これに対し、第1の実施形態による抽出装置10によれば、「はきれいで」、「は悪くない」のようなN-Gramごとに算出するため、単語の組合せで文意が変わる場合でも所望の要約文の抽出が可能となる。すなわち、単語の順序による文意を反映した要約文の抽出が可能となる。

30

【0033】

また、文単位だけで抽出する例では、一文中に複数の文意がある(例えば、図7に示すように一文中にP(ポジティブ)、N(ネガティブ)が複数ある)場合に、所望の要約文を抽出することができない。これに対し、第1の実施形態による抽出装置10は、寄与度の算出で、各単語N-Gramにおける算出された確信度の分散値又は標準偏差値、算出された確信度の最大絶対値、又は、算出された確信度のノルム値と用いる。これにより、一文中に複数の文意があっても適切な要約文の抽出が可能になる。

【0034】

<第2の実施形態>

40

本発明の第2の実施形態による支援装置について、図9を用いて説明する。図9は、第2の実施形態による支援装置の構成を示すブロック図である。支援装置1は、表示装置5、及び、記憶装置6が接続されている。

【0035】

支援装置1は、教師ラベルを付与するテキスト(対象テキスト)、及び、機械学習モデルを用いて、教師ラベルの付与を支援するための要約文を当該テキストから要約文を抽出する機能を有する。さらに、支援装置1は、支援装置1に接続される表示装置5を介して、抽出した要約文を当該支援システムの利用者に提示する機能を有してもよい。具体的には、表示制御部(図示せず)により、抽出された要約文が表示制御される。また、対象テキストは、支援装置1の通信部(図示せず)を介して取得される。

50

【 0 0 3 6 】

記憶装置 6 は、支援装置 1 が取得する、生成する、又は、算出するための各種データを記憶する機能を有する。

【 0 0 3 7 】

表示装置 5 は、支援装置 1 から出力される、教師ラベルを付与するために抽出された要約文の情報を表示する機能を有する。

【 0 0 3 8 】

第 2 の実施形態による支援装置 1 および記憶装置 6 について、図面を用いて詳細に説明する。

【 0 0 3 9 】

支援装置 1 は、抽出装置 10、言語処理部 20、学習部 30、及び、教師ラベル受付部 50 を備える。さらに、抽出装置 10 は、要約文抽出部 40 を備える。なお、第 2 の実施形態の支援装置の説明にあたり、第 1 の実施形態と同じ構成については、同じ符号を付与し、その説明を簡略化する。

【 0 0 4 0 】

支援装置 1 の言語処理部 20 は、教師データを付与するテキストである対象テキストを取得し、取得した対象テキストを単語ごとに分割し、対象テキストを構成する単語、及び、単語区切りのテキストデータを生成する機能を有する。言語処理部 20 は、生成した単語区切りの対象テキストを、抽出装置 10 の要約文抽出部 40 へ渡す、あるいは、記憶装置 6 のテキスト記憶部 62 に記憶させる。

【 0 0 4 1 】

支援装置 1 の学習部 30 は、単語区切りの対象テキストを取得し、単語記憶部 61 に記憶された対象テキストを構成する単語によってインデックス化する。さらに学習部 30 は、単語 N - G r a m ごとの特徴ベクトルを作成後、学習モデル記憶部 63 に格納された学習モデルをパラメータ記憶部 64 から読み込んだパラメータに沿って学習させる。ここでパラメータとは、学習モデルの作成に用いる教師データ（P / N 情報が既知のテキスト、及び、P / N 情報（スコア）等である。なお、単語区切りの対象テキストは、言語処理部 20 から取得してもよく、又は、記憶装置 6 のテキスト記憶部 62 から取得してもよい。

【 0 0 4 2 】

支援装置 1 の教師ラベル受付部 50 は、教師ラベルを付与するテキストに対して、支援システムの利用者によって判定された教師ラベルを受け、テキスト記憶部 62 に判定された教師ラベルの結果を保存する。判定された教師ラベルの受付としては、一般的な入力装置が適用可能である。例えば、マウス、キーボード、又は、タッチパネルなどを用いることができる。

【 0 0 4 3 】

次に、支援装置 1 に接続された記憶装置 6 の構成について図 9 を用いて説明する。記憶装置 6 は、単語記憶部 61、テキスト記憶部 62、学習モデル記憶部 63、及び、パラメータ記憶部 64 を備える。

【 0 0 4 4 】

記憶装置 6 の単語記憶部 61 は、支援装置 1 に入力された対象テキストを構成する単語を記憶する。

【 0 0 4 5 】

記憶装置 6 のテキスト記憶部 62 は、支援装置 1 に入力された対象テキスト又は単語区切りの対象テキストと、対象テキストと対となる教師ラベルと、を記憶する。

【 0 0 4 6 】

記憶装置 6 の学習モデル記憶部 63 は、支援装置 1 に入力された対象テキストを学習するための学習モデルを記憶する。

【 0 0 4 7 】

記憶装置 6 のパラメータ記憶部 64 は、学習モデルの作成と学習に使用するパラメータを記憶する。

10

20

30

40

50

【 0 0 4 8 】

なお、記憶装置 6 が、支援装置 1 の外部に配置され支援装置 1 と接続された例を用いているが、記憶装置 6 が、支援装置 1 の内部に配置され支援装置 1 と接続されていてもよい。

【 0 0 4 9 】

次に、本発明の第 2 の実施形態による支援装置 1 の動作について図面を用いて説明する。図 1 0 は、本発明の第 2 の実施形態による支援装置 1 の動作を示すフローチャートである。

【 0 0 5 0 】

図 1 0 に示すように、支援装置 1 は、教師ラベルを付与するテキスト（対象テキスト）を取得する。支援装置 1 の言語処理部 2 0 は、取得した対象テキストに対し対象テキストを構成する単語で区切った単語区切りの対象テキストを文単位に分割する（S 2 0 1）。図 1 1 は、言語処理（S 2 0 1）のサブルーチンの動作の示すフローチャートである。言語処理部 2 0 は、取得した対象テキストに対して形態素解析を実施して対象テキストを単語区切りに分割する（S 2 0 1 1）。言語処理部 2 0 は、分割した単語、及び、単語区切りの対象テキストをそれぞれ要約文抽出部 4 0 に送る。なお、要約文抽出部 4 0 に送るのではなく、分割した単語、及び、単語区切りの対象テキストをそれぞれ記憶装置（図示せず）に一時的に保存してもよい。

10

【 0 0 5 1 】

図 1 2 は、言語処理（S 2 0 1）のサブルーチンの動作の別の例を示すフローチャートである。図 1 2 に示すサブルーチンの動作では、言語処理部 2 0 は、図 1 1 の形態素解析（S 2 0 1 1）の後に、形態素の係り受けを分析する構文解析を実施する（S 2 0 1 2）。言語処理部 2 0 が、構文解析を実施することで、単語の係り受けの情報が得られ、後段の抽出装置 1 0 における単語 N - G r a m の確信度の算出時に付加的な情報を与えることができ、要約抽出の適切さがより向上することになる。

20

【 0 0 5 2 】

なお、言語処理のステップ（S 2 0 1）では、対象テキストの単語区切りのために形態素解析を用いる例を示したが、単語区切りの対象テキストを生成できるのであれば、形態素解析以外を用いてもよい。また、英語のテキストのように予め分かち書きとなっている対象テキストの場合、対象テキストを文単位で分割する処理をすればよい。

30

【 0 0 5 3 】

分割した単語は、単語記憶部 6 1 に記憶され、複数の文を含むテキストから抽出された文である単語区切りにした対象テキストは、テキスト記憶部 6 2 に記憶される。もしくは、言語処理部 2 0 により、後段の抽出装置 1 0 の要約文抽出部 4 0 へ送られる。

【 0 0 5 4 】

次に、支援装置 1 の抽出装置 1 0 は、文ごとに 1 以上の単語を含む単語集合を生成し、生成された単語集合に対し学習モデルを用いて評価する評価値を算出し、算出された評価値に基づいて文から要約文を抽出する（S 2 0 2）。具体的には、抽出装置 1 0 は、分割された文ごとに N 個の単語をつなげた単語 N - G r a m（N は 2 以上の自然数）を生成し、生成された単語 N - G r a m に対し学習モデルを用いて教師ラベルらしさを表す確信度を算出する。続いて、抽出装置 1 0 は、算出された確信度に基づいて分割された文ごとに教師ラベル判定寄与度を算出し、教師ラベル判定寄与度に応じて要約文を抽出する。また抽出装置 1 0 は、抽出された要約文を表示装置 5 に出力する。

40

【 0 0 5 5 】

抽出装置 1 0 による要約文の抽出のステップは、第 1 の実施形態の抽出装置 1 0 の動作と同様のため、詳細な説明は省略する。なお、要約文の抽出のために、言語処理部 2 0 で生成された単語区切りの対象テキストは、記憶装置 6 のテキスト記憶部 6 2 から取得してもよく、言語処理部 2 0 から取得してもよい。

【 0 0 5 6 】

次に、表示装置 5 は、支援装置 1 の抽出装置 1 0 から出力された要約文を表示する（S

50

203)。図13は、第2の実施形態による支援装置1が表示装置5に出力する画面表示を示す図である。図13に示すように、表示装置5の表示画面は、「テキスト一覧」、「オプション」、「教師ラベル」、「テキスト」の4つの表示エリアで構成されている。

【0057】

「テキスト一覧」の表示エリアは、対象テキスト（教師ラベルを付与するテキスト）を一覧表示する。テキスト一覧で表示する対象テキストは、支援装置1に入力された順でもよく、あるいは所定の降順であってもよい。支援システムの利用者は、「テキスト一覧」の表示エリアに表示された中から対象テキストを選択する。

【0058】

「テキスト」の表示エリアは、「テキスト一覧」で選択されたテキストを表示するエリアである。表示制御部は、抽出装置10から出力された要約文ごとにその文中で、教師単語集合らしいか否かに応じて表示を変化させる。図13中、抽出された要約文ごとに、その文中に「ポジティブ」であると推定した箇所には下線がひかれている。

【0059】

「オプション」の表示エリアは、対象テキストを要約文表示モード、又は、テキスト全文表示モードを選択するエリアであり、これの選択状態によって「テキスト」の表示エリアに表示する内容を変更する。この選択肢は、機械学習が不足している支援システムの初期段階において要約文の抽出が妥当でないときの問題を回避するために用意している。つまり、支援システムの稼働初期はテキスト全文を表示し、教師ラベルが十分に揃うようになれば要約文を表示するよう選択することで、教師ラベルの付与において効率のよい支援システムの運用が可能となる。

【0060】

「教師ラベル」の表示エリアは、「テキスト」に表示されている対象テキストに対して、教師ラベルを設定/変更するためのエリアである。図13中、教師ラベルは、「ポジティブ」に選択されている。

【0061】

支援装置1は、抽出した要約文を表示装置5で表示することで、支援装置1の利用者は、短時間でテキストの内容を理解することができ、教師ラベルの判定の時間と手間を軽減することが可能となる。

【0062】

次に、支援装置1の教師ラベル受付部50は、支援システムの利用者によって判定された教師ラベルの結果を受付ける（S204）。教師ラベル受付部50は、支援システムの利用者が判定した教師ラベルを対象テキストに紐づけてテキスト記憶部62へ記憶させる。

【0063】

続いて、支援装置1の学習部30は、教師ラベル受付部で受付けた教師ラベルと、それに紐づいた対象テキストを用いて、学習モデル記憶部63に記憶された学習モデルを学習させる（S205）。

【0064】

図14は、図10に示す学習のサブルーチンの動作を示すフローチャートである。なお、ここで言う学習は、対象テキストを学習するための特徴ベクトル作成処理を含む（特徴ベクトル作成処理が、学習処理と一体になっているアルゴリズムが存在するため）。

【0065】

まず、学習部30は、学習に用いるベクトルを作成する（S2051）。一般に自然言語の特徴ベクトルは非常に大きいベクトル長のデータであり、そのままでは後段の学習および判別への適用が困難となる。そのため、特徴となる項のみを選択し、圧縮したベクトルを生成する。例えば、特徴ベクトルの生成については、下記論文で詳細に記載されている。

“Sentiment Classification with Supervised Sequence Embedding”, Bespalov, Dmitriy and Qi, Yanjun and Bai, Bing and Shokoufandeh, Ali, Machine Learning and Know

10

20

30

40

50

ledge Discovery in Databases, Vol.7523, pp.159-174, Springer Berlin Heidelberg, 2012, ISBN: 978-3-642-33459-7

上記の論文では、特徴ベクトルの生成を自動で処理する機構を用いている。第2の実施形態では、これに限られず、例えば、主成分分析などにより、重要なベクトル項を分析し、そのベクトル項を選択して、特徴ベクトルを生成する処理をソフトウェアプログラムに組み込んで構成してもよい。

【0066】

続いて、学習部30は、学習モデル記憶部63から学習モデルを読み込み、ベクトル作成のステップ(S2051)によって作成されたベクトルを用いて学習モデルを補正する。学習モデル記憶部63で採用する学習モデルは、任意の教師あり機械学習分類器を適用することでき、この他に、サポートベクタマシン、ニューラルネット、ベイズ分類器などを用いてもよい。

10

【0067】

<第3の実施形態>

本発明の第3の実施形態による支援装置および記憶装置について、図15を用いて説明する。図15は、第3の実施形態による支援装置1及び記憶装置7の構成を示すブロック図である。図15に示すように、第3の実施形態の記憶装置7は、第2の実施形態の記憶装置6と比較して、付加情報記憶部65が追加されている点で相違する。

【0068】

第3の実施形態による記憶装置7の付加情報記憶部65は、対象テキストに関する属性情報を記憶する。これにより、第3の実施形態による支援装置1は、対象テキストの属性情報を使用した学習が可能となる。第2の実施形態の例では、単語記憶部61とテキスト記憶部62で記憶する対象テキストを構成する単語、及び、単語区切り対象テキストのデータを学習部30へ引き渡すことで学習する。それに加えて第3の実施形態の例では、対象テキストのジャンル(論文、小説等)、作者のドメイン(性別、年齢等)、レイアウト(テキスト全体で見た文の出現箇所、文字数)といった付加情報を学習する。これにより、要約文を抽出する精度が向上する。

20

【0069】

<第4の実施形態>

本発明の第4の実施形態である表示制御装置について、図面を用いて説明する。図16は、第4の実施形態による表示制御装置110の構成を示すブロック図である。第5の実施形態の表示制御装置110は、テキストに教師ラベルを付与する者に対して、その教師ラベルの判定を支援するための表示制御装置である。

30

【0070】

図16に示すように、表示制御装置110は、表示制御部140を備える。表示制御部140は、複数の文を含むテキストから学習モデルを用いて算出された評価値に基づいて文から要約文を抽出し、要約文を評価値に基づいた順序で表示制御する。複数の文を含むテキストから学習モデルを用いて算出された評価値に基づいて文から要約文を抽出する点は、第1の実施形態による抽出装置10と同様である。

【0071】

図17は、第4の実施形態による表示制御装置110の動作を示すフローチャートである。図17に示すように、表示制御装置110は、複数の文を含むテキストから学習モデルを用いて算出された評価値に基づいて文から要約文を抽出し(S111)、要約文を評価値に基づいた順序で表示制御する複数の文を含むテキストから文を抽出する(S112)。なお、複数の文を含むテキストから学習モデルを用いて算出された評価値に基づいて文ごとに教師ラベル判定寄与度を算出し、寄与度に応じて要約文を抽出した場合、要約文を寄与度に基づいた順序で表示制御してもよい。

40

【0072】

第4の実施形態によれば、単語の語順に伴う文意を反映した要約文を表示制御することができる。これにより、テキストに教師ラベルを付与する者に対して、その教師ラベルの

50

判定を支援することが可能になる。

【0073】

(ハードウェア構成)

図18は、本発明の第1の実施形態による抽出装置10、第2、3の実施形態による支援装置1、又は第4の実施形態による表示制御装置110をコンピュータ装置で実現したハードウェア構成を示す図である。

【0074】

図18に示すコンピュータ装置は、CPU(Central Processing Unit)91、ネットワーク接続用の通信I/F(通信インターフェース)92、メモリ93、及び、プログラムを格納するハードディスク等の記憶装置94を含む。また、コンピュータ装置は、システムバス97を介して入力装置95及び、出力装置96に接続されている。

【0075】

CPU91は、オペレーティングシステムを動作させて、第1の実施形態による抽出装置10の要約文抽出部40、第2の実施形態による支援装置1の言語処理部20、学習部30、教師ラベル受付部50又は第4の実施形態による表示制御装置の表示制御部140を制御する。またCPU91は、例えば、ドライブ装置に装着された記録媒体からメモリ93にプログラムやデータを読み出す。また、CPU91は、例えば、各実施形態における情報信号を処理する機能を有し、プログラムに基づいて各種機能の処理を実行する。

【0076】

記憶装置94は、例えば、光ディスク、フレキシブルディスク、磁気光ディスク、外付けハードディスク、又は半導体メモリ等である。記憶装置94の一部の記憶媒体は、不揮発性記憶装置であり、そこにプログラムを記憶する。また、プログラムは、通信網に接続されている。図示しない外部コンピュータからダウンロードされてもよい。

【0077】

入力装置95は、例えば、マウス、キーボード、内臓のキーボタン、カード取込口、又は、タッチパネルなどで実現され、入力操作に用いられる。

【0078】

出力装置96は、例えば、ディスプレイで実現され、CPU91により処理された情報等を出力して確認するために用いられる。

【0079】

以上のように、本発明の各実施形態は、図18に示されるハードウェア構成によって実現される。但し、抽出装置10、又は、支援装置1が備える各部の実現手段は、特に限定されない。すなわち、抽出装置10、又は、支援装置1は、物理的に結合した一つの装置により実現されてもよいし、物理的に分離した二つ以上の装置を有線又は無線で接続し、これら複数の装置により実現してもよい。

【0080】

以上、実施形態(及び実施例)を参照して本願発明を説明したが、本願発明は上記実施形態(及び実施例)に限定されるものではない。本願発明の構成や詳細には、本願発明のスコop内で当業者が理解し得る様々な変更をすることができる。

【0081】

上記の実施形態の一部又は全部は、以下の付記のように記載されうるが、以下には限られない。

【0082】

(付記1)

教師ラベルを付与するテキストである対象テキストに対し前記対象テキストを構成する単語で区切った単語区切りの対象テキストを文単位に分割し、前記分割された文ごとにN個の単語をつなげた単語N-Gram(Nは2以上の自然数)を生成し、前記生成された単語N-Gramに対し学習モデルを用いて教師ラベルらしさを表す確信度を算出し、前記算出された確信度に基づいて前記分割された文ごとに教師ラベル判定寄与度を算出し、

10

20

30

40

50

前記寄与度に応じて要約文を抽出する要約文抽出部を備える抽出装置。

【0083】

(付記2)

前記要約文抽出部は、

前記単語 N - G r a m よりも単語区切りが多い単語 M - G r a m ($M > N$: M、N は 2 以上の自然数) の単位、K 文字ごと (K は 1 以上の自然数)、行単位 (改行文字)、ページ単位 (改ページコード)、約物単位、又は、節・段落単位により、前記単語区切りの対象テキストを文単位に分割する、

付記 1 に記載の抽出装置。

【0084】

(付記3)

前記学習モデルは、

スコア情報が既知の教師データであるテキストを用いた単語 N - G r a m が作成され、前記作成された単語 N - G r a m ごとに単語に紐づく特徴ベクトルに置換され、前記特徴ベクトルと対応する前記スコア情報とにより任意の教師あり機械学習分類器に学習させたモデルである、

付記 1 又は付記 2 に記載の抽出装置。

【0085】

(付記4)

前記学習モデルは、任意の教師あり機械学習分類器であり、サポートベクタマシン、ニューラルネットワーク、又は、ベイズ分類器のいずれかである、

付記 1 から 3 のいずれか 1 つに記載の抽出装置。

【0086】

(付記5)

前記寄与度の算出は、各単語 N - G r a m における算出された確信度の分散値又は標準偏差値、各単語 N - G r a m における算出された確信度の最大絶対値、又は、各単語 N - G r a m における算出された確信度のノルム値のいずれかを用いる、

付記 1 から 4 のいずれか 1 つに記載の抽出装置。

【0087】

(付記6)

前記要約文は、前記算出された寄与度が、所定の閾値以上である文、あるいは、前記算出された寄与度を降順に整列したうちの上位数十パーセントとなる文、を抽出する、

付記 1 から 5 のいずれか 1 つに記載の抽出装置。

【0088】

(付記7)

前記抽出装置を含む、

付記 1 ~ 付記 6 のいずれか 1 つに記載の支援装置。

【0089】

(付記8)

言語処理部を備え、

前記言語処理部は、前記単語区切りの対象テキストを生成する、

付記 7 に記載の支援装置。

【0090】

(付記9)

学習部を備え、

前記学習部は、スコア情報が既知の教師データであるテキストを用いた単語 N - G r a m を作成し、前記作成された単語 N - G r a m ごとに単語に紐づく特徴ベクトルに置換し、前記特徴ベクトルと対応する前記スコア情報とにより任意の教師あり機械学習分類器に学習させる、

付記 7 又は付記 8 に記載の支援装置。

10

20

30

40

50

【 0 0 9 1 】

(付記 1 0)

教師ラベル受付部を備え、

前記教師ラベル受付部は、前記対象テキストに対して、前記支援装置の利用者によって判定された教師ラベルを受付ける、

付記 7 ~ 付記 9 のいずれか 1 つに記載の支援装置。

【 0 0 9 2 】

(付記 1 1)

前記支援装置に記憶装置が接続され、

前記記憶装置は、単語記憶部、テキスト記憶部、学習モデル記憶部、及び、パラメータ記憶部を有する、

付記 7 ~ 付記 1 0 のいずれか 1 つに記載の支援装置。

10

【 0 0 9 3 】

(付記 1 2)

前記記憶装置は、付加情報記憶部を有する、

付記 1 1 に記載の支援装置。

【 0 0 9 4 】

(付記 1 3)

前記記憶装置を備える、

付記 1 1 又は付記 1 2 に記載の支援装置。

20

【 0 0 9 5 】

(付記 1 4)

前記支援装置に表示装置が接続され、

前記表示装置は、

付記 7 ~ 付記 1 3 のいずれか 1 つに記載の支援装置。

【 0 0 9 6 】

(付記 1 5)

前記表示装置を備える、

付記 7 ~ 付記 1 4 のいずれか 1 つに記載の支援装置。

【 0 0 9 7 】

(付記 1 6)

教師ラベルを付与するテキストである対象テキストに対し前記対象テキストを構成する単語で区切った単語区切りの対象テキストを文単位に分割し、

前記分割された文ごとに N 個の単語をつなげた単語 N - G r a m (N は 2 以上の自然数) を生成し、

前記生成された単語 N - G r a m に対し学習モデルを用いて教師ラベルらしさを表す確信度を算出し、

前記算出された確信度に基づいて前記分割された文ごとに教師ラベル判定寄与度を算出し、前記寄与度に応じて要約文を抽出する、

抽出方法。

40

【 0 0 9 8 】

(付記 1 7)

コンピュータに、

教師ラベルを付与するテキストである対象テキストに対し前記対象テキストを構成する単語で区切った単語区切りの対象テキストを文単位に分割し、

前記分割された文ごとに N 個の単語をつなげた単語 N - G r a m (N は 2 以上の自然数) を生成し、

前記生成された単語 N - G r a m に対し学習モデルを用いて教師ラベルらしさを表す確信度を算出し、

前記算出された確信度に基づいて前記分割された文ごとに教師ラベル判定寄与度を算出

50

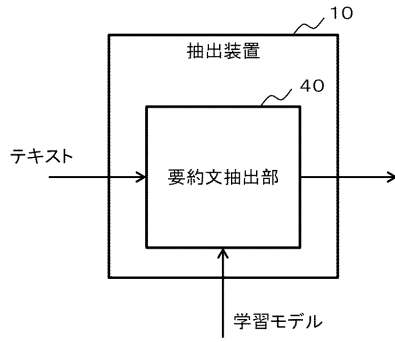
し、前記寄与度に応じて要約文を抽出する、
ことを実行させるための抽出プログラム。

【符号の説明】

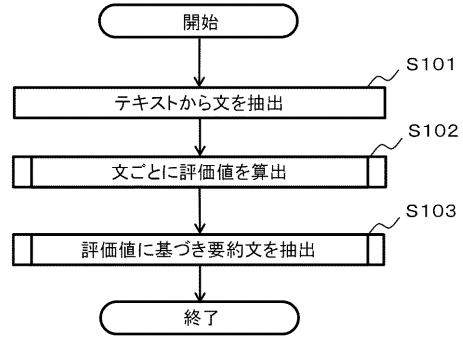
【 0 0 9 9 】

1	支援装置	
5	表示装置	
6	記憶装置	
7	記憶装置	
10	抽出装置	
20	言語処理部	10
30	学習部	
40	要約文抽出部	
50	教師ラベル受付部	
61	単語記憶部	
62	テキスト記憶部	
63	学習モデル記憶部	
64	パラメータ記憶部	
65	付加情報記憶部	
91	CPU	
92	通信I/F(通信インターフェース)	20
93	メモリ	
94	記憶装置	
95	入力装置	
96	出力装置	
97	システムバス	
110	表示制御装置	
140	表示制御部	

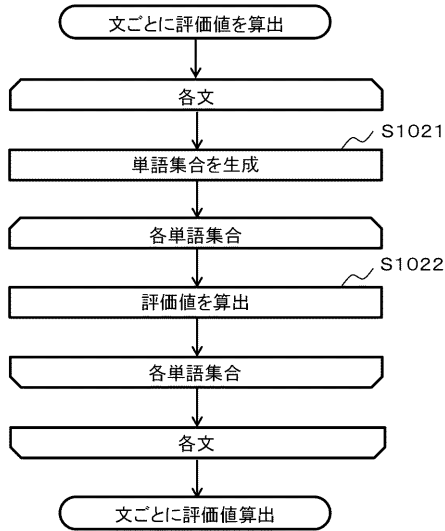
【図1】



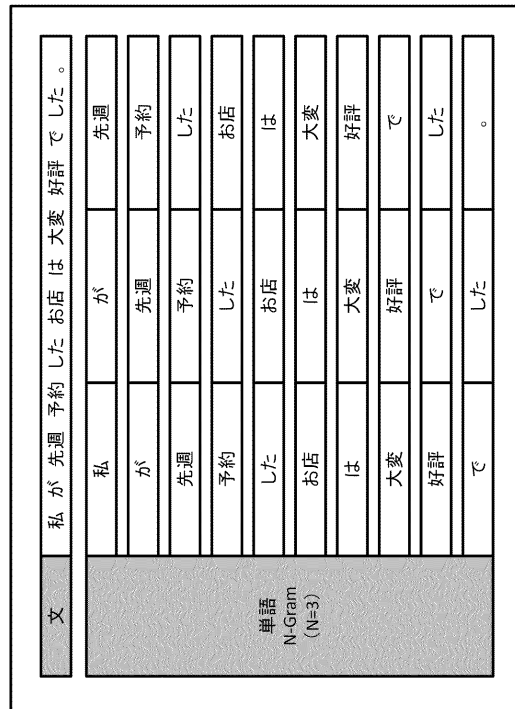
【図2】



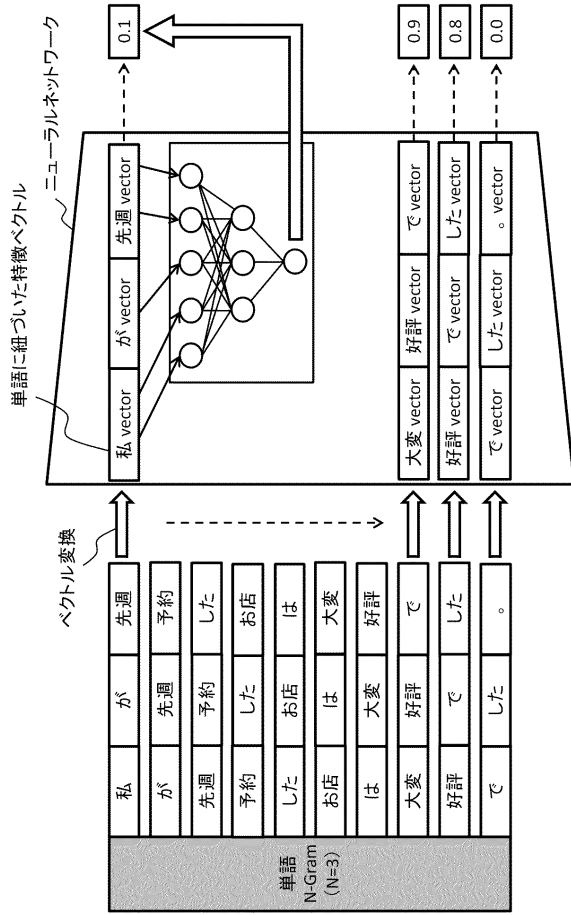
【図3】



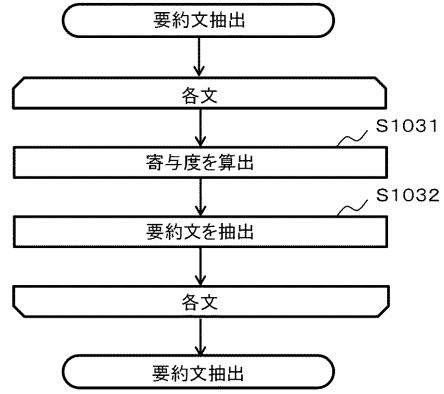
【図4】



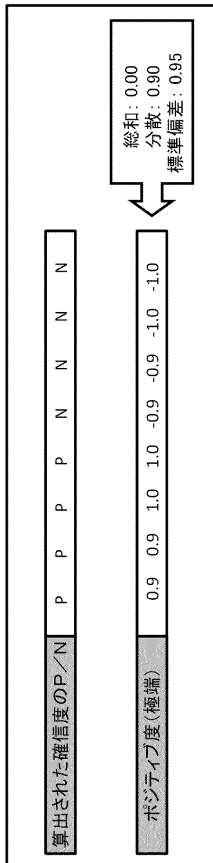
【図5】



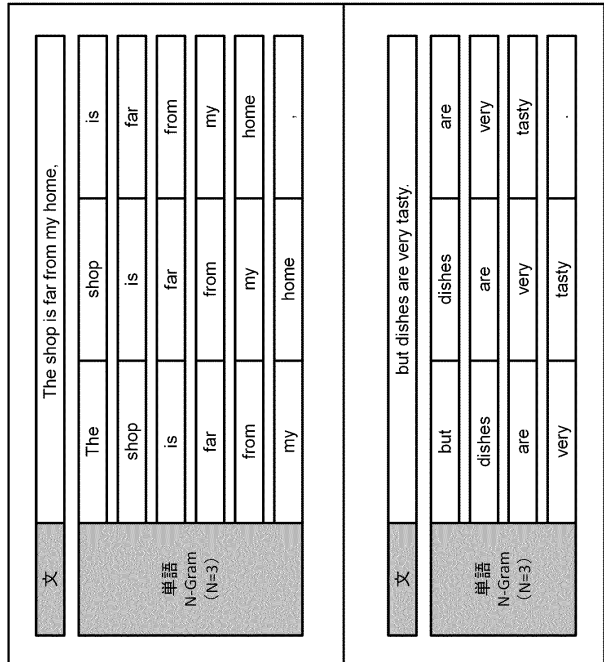
【図6】



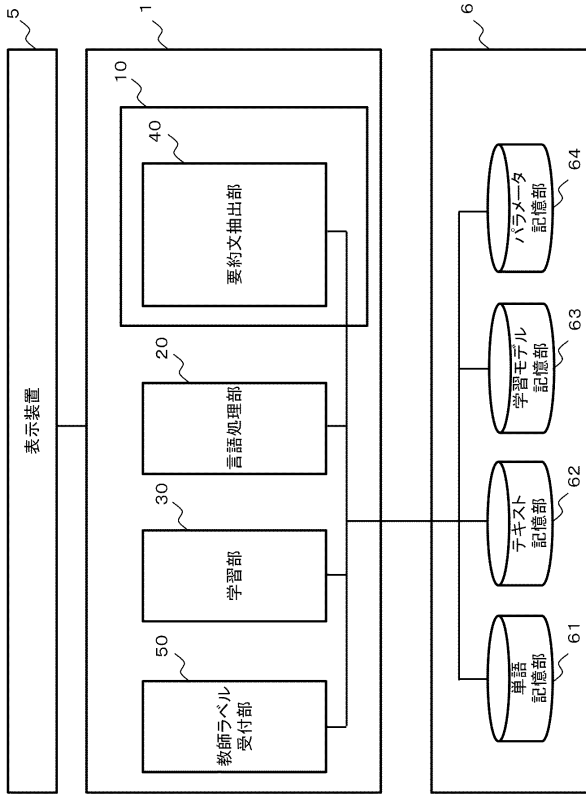
【図7】



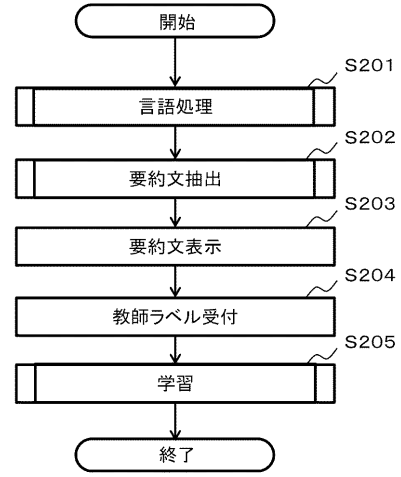
【図8】



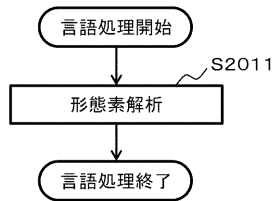
【図9】



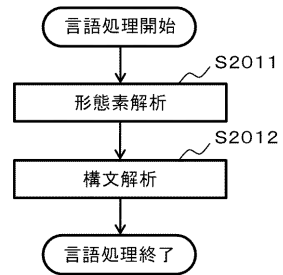
【図10】



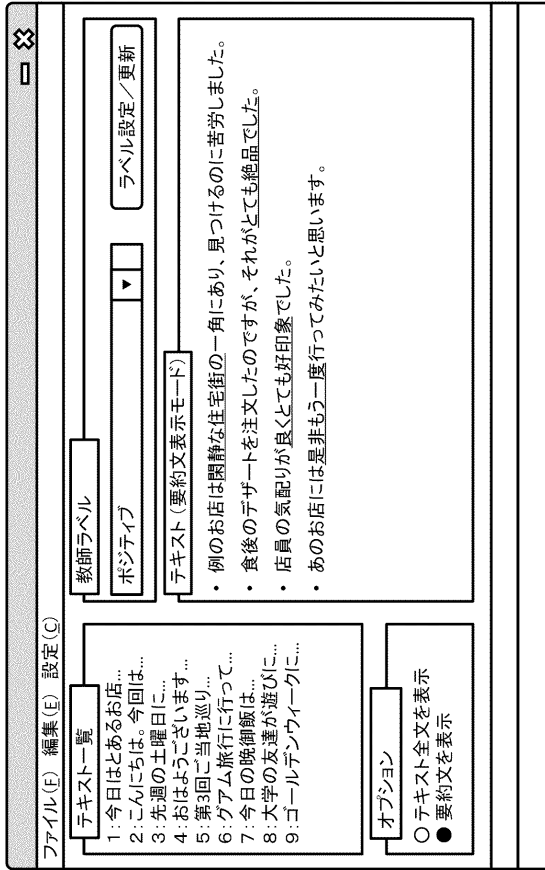
【図11】



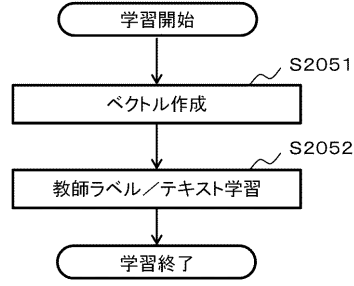
【図12】



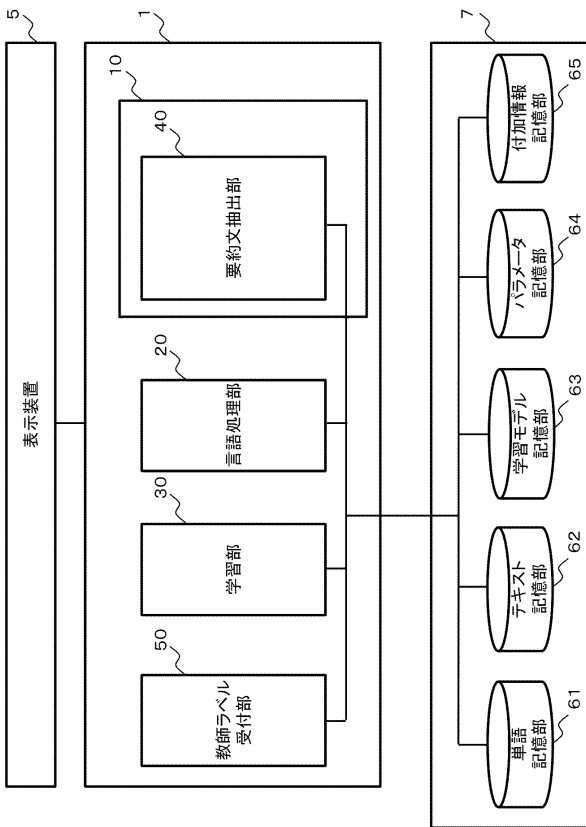
【図 13】



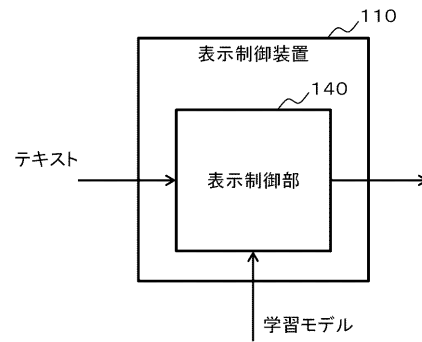
【図 14】



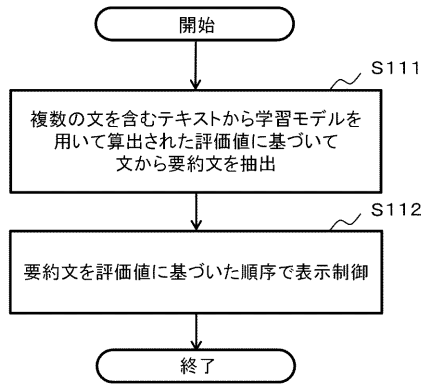
【図 15】



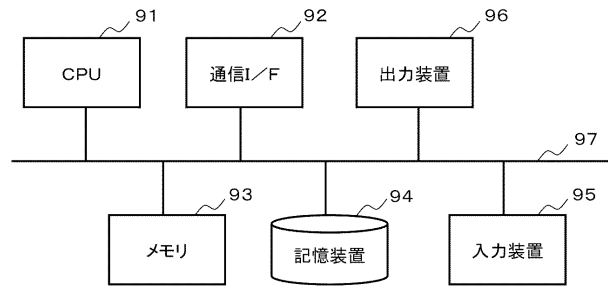
【図 16】



【図17】



【図18】



フロントページの続き

(56)参考文献 特開2003-36262(JP,A)
特開平11-219361(JP,A)
特開2003-337821(JP,A)

(58)調査した分野(Int.Cl., DB名)
G06F16/00-16/958