

(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2013-97723

(P2013-97723A)

(43) 公開日 平成25年5月20日(2013.5.20)

(51) Int.Cl.	F I	テーマコード (参考)
<b>G06F 17/30 (2006.01)</b>	G06F 17/30 220A	5B109
<b>G06F 17/21 (2006.01)</b>	G06F 17/30 170A	
	G06F 17/21 550A	
	G06F 17/30 220Z	

審査請求 未請求 請求項の数 5 O L (全 13 頁)

(21) 出願番号	特願2011-242529 (P2011-242529)	(71) 出願人	000004226 日本電信電話株式会社 東京都千代田区大手町二丁目3番1号
(22) 出願日	平成23年11月4日 (2011.11.4)	(74) 代理人	110001519 特許業務法人太陽国際特許事務所
		(72) 発明者	西川 仁 東京都千代田区大手町二丁目3番1号 日本電信電話株式会社内
		(72) 発明者	牧野 俊朗 東京都千代田区大手町二丁目3番1号 日本電信電話株式会社内
		(72) 発明者	松尾 義博 東京都千代田区大手町二丁目3番1号 日本電信電話株式会社内
		Fターム(参考)	5B109 QA05

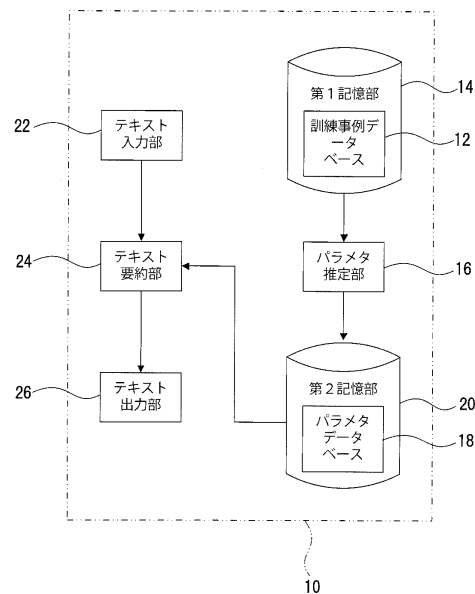
(54) 【発明の名称】 テキスト要約装置、方法及びプログラム

(57) 【要約】

【課題】 ROUGE 等のような要約の品質を評価する評価尺度の値が直接最大化されるようにパラメタを学習する。

【解決手段】 第1記憶部14に記憶された訓練事例データベース12には、テキスト $x_i$ とテキスト $x_i$ から予め生成された要約 $y_i$ とを対応付けた訓練事例が格納されており、パラメタ推定部16は、訓練事例のテキスト $x_i$ から要約 $y'$ を生成し、訓練事例の要約 $y_i$ に対する生成した要約 $y'$ の誤差を、要約の品質を評価する評価尺度を用いて求め、要約 $y$ の生成に用いるパラメタを、求めた誤差が大きくなる程更新の幅が大きくなるように更新し、パラメタデータベース18に格納する。パラメタデータベース18に格納されたパラメタは、テキスト要約部24による要約対象のテキスト $x$ からの要約 $y$ の生成に用いられる。

【選択図】 図1



## 【特許請求の範囲】

## 【請求項 1】

テキストと、当該テキストから予め生成された要約と、を対応付けた訓練事例を格納した訓練事例データベースを記憶する第 1 記憶部と、

前記訓練事例に含まれる前記テキストと対応付けられた要約の品質に対する、要約対象のテキストからの要約の生成に用いられるパラメタの更新前の値を用いて、前記訓練事例に含まれる前記テキストから生成された要約の品質の誤差を、要約の品質を評価する評価尺度を用いて求め、前記パラメタの更新前の値を、求めた前記誤差が大きくなる程更新の幅が大きくなるように更新することで、前記パラメタを推定するパラメタ推定部と、

を含むテキスト要約装置。

10

## 【請求項 2】

前記パラメタ推定部は、前記パラメタとしての重みベクトルの更新前の値を  $w_{old}$ 、前記訓練事例に含まれる前記テキストを  $x_i$ 、前記訓練事例に含まれる前記要約を  $y_i$ 、前記パラメタとしての重みベクトルの更新前の値を用いて前記訓練事例に含まれる前記テキスト  $x_i$  から生成された要約を  $y'$ 、テキスト  $x$  及び要約  $y$  の特徴ベクトルを  $f(x, y)$ 、前記評価尺度を  $ROUGE$ 、前記誤差を  $loss(y'; y_i)$  としたときに、

## 【数 1】

$$w_{new} = w_{old} + \lambda (f(x_i, y_i) - f(x_i, y')) \quad \dots(1)$$

但し、

20

$$\lambda = \frac{loss(y'; y_i) - w_{old} \cdot f(x_i, y_i) + w_{old} \cdot f(x_i, y')}{|f(x_i, y_i) - f(x_i, y')|^2} \quad \dots(2)$$

であり、

$$loss(y'; y_i) = 1 - ROUGE(y'; y_i) \quad \dots(3)$$

上記(1)~(3)式に従って前記パラメタとしての重みベクトルの更新後の値  $w_{new}$  を演算することで前記パラメタを推定する請求項 1 記載のテキスト要約装置。

## 【請求項 3】

前記パラメタ推定部によって推定された前記パラメタを格納するパラメタ・データベースを記憶する第 2 記憶部と、

30

要約対象のテキストを受け付けるテキスト入力部と、

前記テキスト入力部によって受け付けられた前記要約対象のテキストと、前記パラメタ・データベースに格納された前記パラメタと、に基づいて、前記要約対象のテキストの要約を生成するテキスト要約部と、

前記テキスト要約部によって生成された要約をテキストとして出力するテキスト出力部と、

を更に備えた請求項 1 又は請求項 2 記載のテキスト要約装置。

## 【請求項 4】

テキストと、当該テキストから予め生成された要約と、を対応付けた訓練事例を格納した訓練事例データベースが第 1 記憶部に記憶された状態で、

40

前記訓練事例に含まれる前記テキストと対応付けられた要約の品質に対する、要約対象のテキストからの要約の生成に用いられるパラメタの更新前の値を用いて、前記訓練事例に含まれる前記テキストから生成された要約の品質の誤差を、要約の品質を評価する評価尺度を用いて求め、前記パラメタの更新前の値を、求めた前記誤差が大きくなる程更新の幅が大きくなるように更新することで、前記パラメタを推定するパラメタ推定ステップ

を含むテキスト要約方法。

## 【請求項 5】

コンピュータを、請求項 1 ~ 請求項 3 の何れか 1 項記載のテキスト要約装置を構成する各手段として機能させるためのテキスト要約プログラム。

50

## 【発明の詳細な説明】

## 【技術分野】

## 【0001】

本発明はテキスト要約装置、テキスト要約方法及びテキスト要約プログラムに関する。

## 【背景技術】

## 【0002】

近年、電子化されたテキストが大量に流通するようになってきたことを背景として、それらのテキストに記述されている情報を迅速に把握するために、コンピュータ等の機械によってテキストを要約させる(要約を生成させる)技術に対するニーズが高まっている。

## 【0003】

機械によってテキストを要約させる技術では、一般に、要約対象のテキストから、その内容を代表していると思われる文(重要文)、或いは、そのような文の集合(重要文集合)を1つ以上選び出し、それらを連結することで要約が生成される。文、或いは文の集合に対して要約対象のテキストの内容を代表しているか否かを評価する際には、各々の文、或いは文の集合を特徴ベクトルとして表現し、この特徴ベクトルと、予め何らかの方法で推定した重みベクトル(以下、パラメタともいう)と、の内積がスコアとして算出される(例えば非特許文献1も参照)。

## 【0004】

また、文、或いは文の集合を表現した特徴ベクトルとの内積を求めるパラメタは、要約対象のテキスト、或いは要約対象のテキストの集合と、それらに対応する要約と、から成るペアの集合(以下、訓練事例という)に基づいて予め学習される。パラメタの学習に際し、パラメタは何らかの誤差関数を最小化するように学習される。

## 【先行技術文献】

## 【非特許文献】

## 【0005】

【非特許文献1】Wen-tau Yih, Joshua Goodman, Lucy Vanderwende and Hisami Suzuki, "Multi-Document Summarization by Maximizing Informative Content-Words.", In Proceedings of International Joint Conference on Artificial Intelligence(IJCAI), 2007.

【非特許文献2】Koby Crammer, Ofel Dekel, Joseph Keshet, Shai Shalev-Shwartz and Yooram Singer, "Online Passive-Aggressive Algorithms.", Journal of Machine Learning Research, Vol.7, 2006.

## 【発明の概要】

## 【発明が解決しようとする課題】

## 【0006】

しかしながら、テキストを要約する技術で用いるパラメタの学習に際し、例えば非特許文献1に記載の技術では、要約の品質が直接最大化されるようなパラメタの学習を行っておらず、要約に含まれ易い文が高いスコアになるようにパラメタの学習を行っている。機械によって生成された要約の品質を評価するための尺度としては、例えばROUGE(非特許文献2を参照)と称される評価尺度が多用されているが、要約に含まれ易い文が高いスコアになるようにパラメタを学習することと、学習したパラメタを用いて生成した要約のROUGE等の評価尺度が高い値になることには直接関係がなく、パラメタの学習精度に関して改善の余地があった。

## 【0007】

本発明は上記事実を考慮して成されたもので、ROUGE等のような要約の品質を評価する評価尺度の値が直接最大化されるようにパラメタを学習できるテキスト要約装置、テキスト要約方法及びテキスト要約プログラムを得ることが目的である。

## 【課題を解決するための手段】

## 【0008】

本発明に係るテキスト要約装置は、テキストと、当該テキストから予め生成された要約

10

20

30

40

50

と、を対応付けた訓練事例を格納した訓練事例データベースを記憶する第1記憶部と、前記訓練事例に含まれる前記テキストと対応付けられた要約の品質に対する、要約対象のテキストからの要約の生成に用いられるパラメタの更新前の値を用いて、前記訓練事例に含まれる前記テキストから生成された要約の品質の誤差を、要約の品質を評価する評価尺度を用いて求め、前記パラメタの更新前の値を、求めた前記誤差が大きくなる程更新の幅が大きくなるように更新することで、前記パラメタを推定するパラメタ推定部と、を含んで構成されている。

【0009】

また、本発明において、前記パラメタ推定部は、前記パラメタとしての重みベクトルの更新前の値を  $w_{old}$ 、前記訓練事例に含まれる前記テキストを  $x_i$ 、前記訓練事例に含まれる前記要約を  $y_i$ 、前記パラメタとしての重みベクトルの更新前の値を用いて前記訓練事例に含まれる前記テキスト  $x_i$  から生成された要約を  $y'$ 、テキスト  $x$  及び要約  $y$  の特徴ベクトルを  $f(x, y)$ 、前記評価尺度を  $ROUGE$ 、前記誤差を  $loss(y'; y_i)$  としたときに、

【0010】

【数1】

$$w_{new} = w_{old} + \lambda (f(x_i, y_i) - f(x_i, y')) \quad \dots(1)$$

但し、

$$\lambda = \frac{loss(y'; y_i) - w_{old} \cdot f(x_i, y_i) + w_{old} \cdot f(x_i, y')}{|f(x_i, y_i) - f(x_i, y')|^2} \quad \dots(2)$$

【0011】

であり、

$$loss(y'; y_i) = 1 - ROUGE(y'; y_i) \quad \dots(3)$$

【0012】

上記(1)~(3)式に従って前記パラメタとしての重みベクトルの更新後の値  $w_{new}$  を演算することで前記パラメタを推定することができる。

【0013】

また、本発明は、前記パラメタ推定部によって推定された前記パラメタを格納するパラメタ・データベースを記憶する第2記憶部と、要約対象のテキストを受け付けるテキスト入力部と、前記テキスト入力部によって受け付けられた前記要約対象のテキストと、前記パラメタ・データベースに格納された前記パラメタと、に基づいて、前記要約対象のテキストの要約を生成するテキスト要約部と、前記テキスト要約部によって生成された要約をテキストとして出力するテキスト出力部と、を更に備えることができる。

【0014】

本発明に係るテキスト要約方法は、テキストと、当該テキストから予め生成された要約と、を対応付けた訓練事例を格納した訓練事例データベースが第1記憶部に記憶された状態で、前記訓練事例に含まれる前記テキストと対応付けられた要約の品質に対する、要約対象のテキストからの要約の生成に用いられるパラメタの更新前の値を用いて、前記訓練事例に含まれる前記テキストから生成された要約の品質の誤差を、要約の品質を評価する評価尺度を用いて求め、前記パラメタの更新前の値を、求めた前記誤差が大きくなる程更新の幅が大きくなるように更新することで、前記パラメタを推定するパラメタ推定ステップを含んでいる。

【0015】

本発明に係るテキスト要約プログラムは、コンピュータを、請求項1~請求項3の何れか1項記載のテキスト要約装置を構成する各手段として機能させる。

【発明の効果】

【0016】

本発明は、 $ROUGE$ 等のような要約の品質を評価する評価尺度の値が直接最大化され

るようにパラメタを学習できる、という効果を有する。

【図面の簡単な説明】

【0017】

【図1】実施形態で説明したテキスト要約装置を示すブロック図である。

【図2】テキスト要約装置として機能するコンピュータの概略ブロック図である。

【図3】テキスト要約処理の概略を示すフローチャートである。

【図4】訓練事例データベースに格納される訓練事例の一例を示す図表である。

【図5】(A)はテキスト及び要約の一例を示す図表、(B)は(A)に示すテキスト及び要約に対応する特徴ベクトルの一例を示す説明図である。

【図6】パラメタ学習処理のアルゴリズムの一例を示すフローチャートである。

【図7】パラメタ・データベースの一例を示す図表である。

【図8】要約対象のテキストの入力形式の一例を示す図表である。

【図9】要約生成処理のアルゴリズムの一例を示すフローチャートである。

【発明を実施するための形態】

【0018】

以下、図面を参照して本発明の実施形態の一例を詳細に説明する。図1には本実施形態に係るテキスト要約装置10が示されている。テキスト要約装置10は、パラメタ(本実施形態では重みベクトル $w$ )を推定するために必要となる訓練事例を格納した訓練事例データベース12を記憶する第1記憶部14と、訓練事例データベース12に格納された訓練事例を受け付けてパラメタを推定するパラメタ推定部16と、パラメタ推定部16で推定されたパラメタを格納したパラメタ・データベース18を記憶する第2記憶部20と、要約の対象とするテキストを受け付けるテキスト入力部22と、テキスト入力部22で受け付けられた要約対象のテキストとパラメタ・データベース18に格納されたパラメタを入力として要約を生成するテキスト要約部24と、テキスト要約部24によって生成された要約を出力するテキスト出力部26と、を備えている。

【0019】

テキスト要約装置10は、例えば図2に示すコンピュータ30で実現することができる。コンピュータ30はCPU32、メモリ34、不揮発性の記憶部36、キーボード38、マウス40、ディスプレイ42を備え、これらはバス44を介して互いに接続されている。なお、記憶部36はHDD(Hard Disk Drive)やフラッシュメモリ等によって実現できる。記録媒体としての記憶部36には、コンピュータ30をテキスト要約装置10として機能させるためのテキスト要約プログラム46、訓練事例データベース12、パラメタ・データベース18が各々記憶されている。CPU32は、テキスト要約プログラム46を記憶部36から読み出してメモリ34に展開し、テキスト要約プログラム46が有するプロセスを順次実行する。

【0020】

テキスト要約プログラム46は、テキスト入力プロセス50、パラメタ推定プロセス52、テキスト要約プロセス54及びテキスト出力プロセス56を有する。CPU32は、テキスト入力プロセス50を実行することで、図1に示すテキスト入力部22として動作する。またCPU32は、パラメタ推定プロセス52を実行することで、図1に示すパラメタ推定部16として動作する。またCPU32は、テキスト出力プロセス56を実行することで、図1に示すテキスト出力部26として動作する。なお、テキスト要約プログラム46は本発明に係るテキスト要約プログラムの一例である。

【0021】

テキスト要約装置10がコンピュータ30で実現される場合、訓練事例データベース12を記憶する記憶部36は第1記憶部14として用いられ、メモリ34の一部領域は訓練事例データベース12として用いられる。また、パラメタ・データベース18を記憶する記憶部36は第2記憶部20として用いられ、メモリ34の一部領域はパラメタ・データベース18として用いられる。これにより、テキスト要約プログラム46を実行したコンピュータ30が、テキスト要約装置10として機能することになる。

## 【 0 0 2 2 】

次に本実施形態の作用として、まず図3を参照し、コンピュータ30のCPU32でテキスト要約プログラム46が実行されることで実現されるテキスト要約処理の概略を説明する。テキスト要約処理のステップ60において、パラメタ推定部16は、訓練事例データベース12から訓練事例を受け付け、受け付けた訓練事例に基づいてパラメタの学習を行う(パラメタ学習の詳細は後述)。また、次のステップ62において、パラメタ推定部16は、学習したパラメタをパラメタ・データベース18に格納する。

## 【 0 0 2 3 】

次のステップ64において、テキスト入力部22は、入力された要約対象のテキストを受け付ける。なお、要約対象のテキストは、例えばキーボード38を介して入力されたテキストでもよいし、例えば記憶部36に予め記憶されたテキストでもよいし、例えば通信回線を介してコンピュータ30と接続された他の機器から受信したテキストでもよい。

## 【 0 0 2 4 】

また、ステップ66において、テキスト要約部24は、テキスト入力部22によって受け付けられた要約対象のテキストと、パラメタ・データベース18に格納されたパラメタと、から要約を生成する(要約生成の詳細は後述)。

## 【 0 0 2 5 】

そしてステップ68において、テキスト出力部26は、テキスト要約部24によって生成された要約を外部へ出力する。要約の出力は、例えばディスプレイ42にテキストとして表示してもよいし、要約を読み上げる音声として出力してもよいし、テキストデータとして記録メディアに記録するか他の機器へ送信するようにしてもよい。

## 【 0 0 2 6 】

また、図3では、パラメタの学習及びパラメタ・データベース18への格納(ステップ60,62)と、要約の生成・出力(ステップ64~68)と、を一連の処理として示しているが、パラメタの学習及びパラメタ・データベース18への格納を行った後、時間を空けて要約の生成・出力を行ってもよい。

## 【 0 0 2 7 】

次に、テキスト要約処理の詳細を説明するにあたり、まず、訓練事例データベース12に格納される訓練事例について、図4を参照して説明する。図4には、訓練事例データベース12に格納される訓練事例の一例が示されている。図4において、訓練事例データベース12は、各レコードに1つの訓練事例を各々格納している。訓練事例データベース12の各レコードに格納されているそれぞれの訓練事例は、個々の訓練事例を識別するための事例番号、要約対象のテキスト、及び、当該要約対象のテキストから生成された要約の各情報を含んでいる。なお、訓練事例に含まれる要約は、例えば人手によって作成され、後述するパラメタ(重みベクトル $w$ )の学習における基準として用いることが可能な正しい要約である。

## 【 0 0 2 8 】

続いて、パラメタ(重みベクトル $w$ )を用いた要約の生成について説明する。要約対象のテキストを $x$ 、要約対象のテキスト $x$ から生成された要約を $y$ とする。要約対象のテキスト $x$ が与えられると、要約対象のテキスト $x$ と、要約対象のテキスト $x$ から生成された要約 $y$ と、から特徴ベクトル $f(x, y)$ が得られる。特徴ベクトル $f(x, y)$ と、パラメタ(重みベクトル $w$ )の内積 $w \cdot f(x, y)$ が要約 $y$ のスコアである。要約を生成する際には、要約対象のテキスト $x$ から生成される要約 $y$ のうち、パラメタ(重みベクトル $w$ )と特徴ベクトル $f(x, y)$ の内積 $w \cdot f(x, y)$ が最大となる要約 $y$ を要約 $y'$ として生成する。上記を式で表すと次の(4)式のようになる。

## 【 0 0 2 9 】

$$y' = \arg y \max w \cdot f(x, y) \dots (4)$$

要約の生成方法の一例については後述するが、任意の要約生成方法を適用することができる。

## 【 0 0 3 0 】

10

20

30

40

50

要約対象のテキスト  $x$  と当該テキスト  $x$  から生成された要約  $y$  から得られる特徴ベクトル  $f(x, y)$  は、各要素が互いに異なる単語に対応し、一例として、要約  $y$  が或る単語を含んでいれば対応する要素の値が 1、含んでいなければ対応する要素の値が 0 となるベクトルを適用することができる。図 5 (A) に示すテキスト  $x$  及び要約  $y$  から得られる特徴ベクトル  $f(x, y)$  の一例を図 5 (B) に示す。図 5 (A) に示す要約  $y$  は「路上」や「強盗」といった単語を含むため、図 5 (B) に示す特徴ベクトル  $f(x, y)$  では、これらの単語に対応する要素の値が 1 となっている。一方、図 5 (A) に示す要約  $y$  は「道路」や「泥棒」といった単語は含んでいないため、図 5 (B) に示す特徴ベクトル  $f(x, y)$  では、これらの単語に対応する要素の値が 0 となっている。

【0031】

10

次に、パラメタ推定部 16 によって行われるパラメタの学習を説明する。  $n$  個の訓練事例が与えられたと仮定し、そのうちの  $i$  番目の訓練事例のテキストを  $x_i$ 、要約を  $y_i$  とする。例えば事例番号 1 のテキストを  $x_1$ 、要約を  $y_1$  とする。

【0032】

パラメタの学習に際しては、  $n$  個の訓練事例が 1 つずつ取り上げられる。その時点でのパラメタ  $w$  の下で、テキスト  $x_i$  から生成された要約  $y'$  と、人手によって作成された正しい要約  $y_i$  が得られた際に、要約  $y'$  の要約としての良さ(品質)は、要約  $y'$  を評価するための評価尺度、一例として ROUGE を用い、要約  $y_i$  を基準とした要約  $y'$  の ROUGE の値 ( $\text{ROUGE}(y'; y_i)$ ) として算出できる。 ROUGE の値は 0 ~ 1 の範囲であるので、前出の (3) 式のように 1 から ROUGE の値を減算することで、要約  $y_i$  を基準(正しい要約)としたときの要約  $y'$  の悪さ、すなわち要約  $y_i$  に対する要約  $y'$  の誤差が得られる。なお、前出の (3) 式では誤差として誤差関数  $\text{loss}(y'; y_i)$  を用いている。

20

【0033】

上記の誤差を減らすようにパラメタを更新できれば、より良い要約を生成できるパラメタを得ることができる。そこで、以下の (5) 式のようにパラメタを更新する。

【0034】

【数 2】

$$w_{new} = \text{arg min}_w \frac{1}{2} |w - w_{old}|^2 \quad \dots (5)$$

但し、  $w_{new} \cdot f(x_i, y_i) - w_{new} \cdot f(x_i, y') \geq \text{loss}(y'; y_i)$

30

【0035】

(5) 式において、  $w_{old}$  は現在(更新前)のパラメタであり、  $w_{new}$  は更新後のパラメタである。上記の (5) 式は、パラメタ  $w$  の変化を最小にしつつ、更新後のパラメタ  $w_{new}$  において、正しい要約  $y_i$  のスコアが現在(更新前)のパラメタ  $w_{old}$  を用いて生成された要約  $y'$  のスコアより必ず大きくなるように、パラメタ  $w$  を更新している。その際、要約  $y_i$  のスコアから要約  $y'$  のスコアを減算した残差が誤差関数  $\text{loss}(y'; y_i)$  以上となるようにしている。上記の (5) 式を解くと、前出のパラメタ更新式((1), (2)式)が得られる。

【0036】

40

前出の更新式((1), (2)式)では、現在(更新前)のパラメタ  $w_{new}$  を、誤差関数に ROUGE を用いて求めた要約  $y_i$  に対する要約  $y'$  の誤差が大きくなる程更新の幅が大きくなるように更新しているため、パラメタ  $w$  の更新に伴って ROUGE の値が直接高くなるようにパラメタ  $w$  を更新(学習)することができる。

【0037】

次に図 6 を参照し、パラメタ推定部 16 によるパラメタ学習処理(図 3 に示すテキスト要約処理のステップ 60, 62 の詳細)の具体的なアルゴリズムの一例を説明する。パラメタ学習処理のステップ 70 において、パラメタ推定部 16 は、事例番号  $i$  が付与された訓練事例に含まれるテキスト  $x_i$  及び要約  $y_i$  を訓練事例データベース 12 から  $n$  個(事例番号  $i = 1 \sim n$  の訓練事例を)取得する。またステップ 72 において、パラメタ推定部 1

50

6は、予め設定された繰り返し回数Tを取得する。次のステップ74において、パラメタ推定部16は、パラメタ(重みベクトル $w$ )を、全ての要素の値が0のベクトルへ初期化する。またステップ76において、パラメタ推定部16は変数 $t$ に1を設定し、次のステップ78において、パラメタ推定部16は変数 $i$ に1を設定する。

【0038】

次のステップ80において、パラメタ推定部16は、事例番号 $i$ の訓練事例について、現在のパラメタ(重みベクトル $w$ )と特徴ベクトル $f(x_i, y)$ との内積 $w \cdot f(x_i, y)$ が最大となる要約 $y$ を、要約 $y'$ として生成する(次の(6)式も参照)。

【0039】

$$y' = \arg y \max w \cdot f(x_i, y) \dots (6)$$

10

【0040】

次のステップ82において、パラメタ推定部16は、まず前出の(3)式により要約 $y_i$ に対する要約 $y'$ の誤差 $\text{loss}(y'; y_i)$ を演算した後に、前出のパラメタ更新式((1),(2)式)によりパラメタ(重みベクトル $w$ )を更新する。これにより、要約 $y_i$ に対する要約 $y'$ の誤差が大きくなる程更新の幅が大きくなるようにパラメタ(重みベクトル $w$ )が更新される。

【0041】

ステップ84において、パラメタ推定部16は、変数 $i$ が先のステップ70で取得した訓練事例の総数 $n$ 以上になったか否か判定する。ステップ84の判定が否定された場合はステップ86へ移行し、ステップ86において、パラメタ推定部16は変数 $i$ を1だけインクリメントした後にステップ80に戻る。これにより、ステップ84の判定が肯定される迄ステップ80~ステップ86が繰り返され、 $n$ 個の訓練事例を1つずつ用いてパラメタ(重みベクトル $w$ )の更新(学習)が $n$ 回繰り返される。

20

【0042】

ステップ84の判定が肯定されるとステップ88へ移行する。ステップ88において、パラメタ推定部16は、変数 $t$ が先のステップ72で取得した繰り返し回数T以上になったか否か判定する。ステップ88の判定が否定された場合はステップ90へ移行し、ステップ90において、パラメタ推定部16は変数 $t$ を1だけインクリメントした後にステップ78に戻る。これにより、ステップ88の判定が肯定される迄ステップ78~ステップ90が繰り返され、 $n$ 個の訓練事例が各々T回ずつ用いられてパラメタ(重みベクトル $w$ )の更新(学習)が繰り返される。これにより、要約の品質を評価する評価尺度であるROUGEの値が直接最大化されるようにパラメタ(重みベクトル $w$ )が更新(学習)される。

30

【0043】

ステップ88の判定が肯定されると、パラメタ(重みベクトル $w$ )の更新を終了してステップ92へ移行し、ステップ92において、パラメタ推定部16は、上記処理で更新(学習)されたパラメタ(重みベクトル $w$ )をパラメタ・データベース18に格納し、パラメタ学習処理を終了する。パラメタ・データベース18の一例を図7に示す。パラメタ・データベース18は、各レコードに、或る単語を含んでいるか否か等の「特徴」と、当該特徴の「重み」の各情報が設定されている。図7に示すパラメタ・データベース18における最初のレコードは、単語「道路」を含んでいるという特徴の重みが-0.03であるということを示している。

40

【0044】

上述したパラメタ学習処理が行われた後、図3に示すテキスト要約処理のステップ64において、テキスト入力部22は、入力された要約対象のテキスト $x$ を受け付ける。テキスト入力部22によって取得されるテキスト $x$ の入力形式の一例を図8に示す。図8に示すテキスト $x$ の入力形式では、テキスト $x$ が一文毎に分割され、更に形態素解析が行われた結果が入力される。図8において、表の各行は各々単一の形態素に対応しており、個々の形態素毎に「表記」「品詞」「読み」及び「標準形」の各情報が付加されている。

【0045】

続いて図9を参照し、テキスト要約部24による要約生成処理(図3に示すテキスト要

50



約処理のステップ66の詳細)の具体的なアルゴリズムの一例を説明する。要約生成処理のステップ100において、テキスト要約部24は、要約対象のテキストx及び要約長Kをテキスト入力部22から取得する。またステップ102において、テキスト要約部24は、パラメタ・データベース18からパラメタ(重みベクトル)wを取得する。また、ステップ104において、テキスト要約部24は、ステップ100で取得したテキストxを文の集合Uに格納し、次のステップ106において、テキスト要約部24は、要約を表す集合Sを空集合に初期化する。

【0046】

ステップ108において、テキスト要約部24は、文の集合Uが空集合よりも大きい(文の集合Uに何らかの文が格納されている)か否かを判定する。ステップ108の判定が肯定された場合はステップ110へ移行する。ステップ110において、テキスト要約部24は、文の集合Uに格納されている文の中から、要約の集合Sに文 $s_i$ を入れた時のスコア(=特徴ベクトルfとパラメタwの内積)と要約の集合Sに文 $s_i$ を入れていない時のスコアの差を、文 $s_i$ の長さで除した値が最大の文 $s_i$ (次の(7)式を満たす文 $s_i$ )を選択する。

【0047】

$$s_i = \arg \max_{s_i \in U} ((w \cdot f(x, \{S, s_i\}) - w \cdot f(x, S)) / \text{length}(s_i)) \quad \dots(7)$$

なお、(7)式において、 $\text{length}(s_i)$ は文 $s_i$ の長さである。

【0048】

次のステップ112において、テキスト要約部24は、要約の集合Sに既に入っている文の長さに、ステップ110で選択した文 $s_i$ の長さを加えた長さが要約長K以下か否かを判定する(次の(8)式も参照)。

$$\text{length}(\{S, s_i\}) \leq K \quad \dots(8)$$

【0049】

ステップ112の判定が肯定された場合はステップ114へ移行する。ステップ114において、テキスト要約部24は、ステップ110で選択した文 $s_i$ を要約の集合Sに加えた後に(次の(9)式も参照)、ステップ116へ移行する。

$$S = \{S, s_i\} \quad \dots(9)$$

【0050】

また、ステップ112の判定が否定された場合は、ステップ114をスキップして(文 $s_i$ を要約の集合Sに加えることなく)ステップ116へ移行する。

【0051】

ステップ116において、テキスト要約部24は、ステップ110で選択した文 $s_i$ を文の集合Uから除去する( $U = U - s_i$ )。ステップ116の処理を行うとステップ108に戻り、ステップ108の判定が肯定される迄ステップ108～ステップ116を繰り返す。文の集合Uが空集合になると、ステップ108の判定が否定されてステップ118へ移行する。

【0052】

ステップ118において、テキスト要約部24は、テキストxの文の中で、長さが要約長K以内で、単一の文として最もスコアの高い文 $s_i$ を文vとして選択する(次の(10)式も参照)。

$$v = \arg \max_{s_i \in x : \text{length}(s_i) \leq K} w \cdot f(x, s_i) \quad \dots(10)$$

【0053】

次のステップ120において、テキスト要約部24は、要約の集合Sに入っている文のスコアが、先のステップ118で選択した文vのスコア以上か否かを判定する。ステップ120の判定が肯定された場合はステップ122へ移行する。ステップ122において、テキスト要約部24は、要約の集合Sに入っている文をテキストxの要約として出力し、要約生成処理を終了する。また、ステップ120の判定が否定された場合はステップ124へ移行する。ステップ124において、テキスト要約部24は、先のステップ118で選

10

20

30

40

50

択した文  $v$  をテキスト  $x$  の要約として出力し、要約生成処理を終了する。

【 0 0 5 4 】

上述した要約生成処理のステップ 1 2 2 又はステップ 1 2 4 において、テキスト要約部 2 4 によって出力されたテキスト  $x$  の要約は、図 3 に示すテキスト要約処理のステップ 6 8 において、テキスト出力部 2 6 により外部へ出力される。

【 0 0 5 5 】

なお、図 9 に要約生成処理として示した要約の生成方法は一例であり、本発明におけるパラメタの更新方法は、任意の要約生成方法と組み合わせることが可能である。

【 0 0 5 6 】

また、上記では、要約の品質を評価する評価尺度として R O U G E を用いた態様を説明したが、本発明はこれに限定されるものではなく、R O U G E 以外の評価尺度を適用することも可能である。

10

【 0 0 5 7 】

また、上記では訓練事例データベース 1 2 が、コンピュータ 3 0 に設けられた記憶部 3 6 に記憶されている態様を説明したが、本発明はこれに限定されるものではなく、通信回線を介してコンピュータ 3 0 と接続された別のコンピュータに設けられた記憶部に記憶されていてもよい。この場合、パラメタの学習にあたり、テキスト  $x_i$  及び要約  $y_i$  を各々含む  $n$  個の訓練事例は、通信回線を介して前記別のコンピュータから受信するように構成することができる。また、要約対象のテキスト  $x$  についても、通信回線を介して前記別のコンピュータから受信する構成であってもよい。

20

【 0 0 5 8 】

更に、上記では、コンピュータ 3 0 がテキスト要約プログラム 4 6 を実行することで、コンピュータ 3 0 がテキスト要約装置 1 0 として機能する態様を説明したが、本発明はこれに限定されるものではなく、図 1 に示した各機能ブロック(パラメタ推定部 1 6、テキスト入力部 2 2、テキスト要約部 2 4 及びテキスト出力部 2 6)は、それぞれハードウェアで構成することも可能である。

【 0 0 5 9 】

また、上記ではテキスト要約プログラム 4 6 が記憶部 3 6 に予め記憶(インストール)されている態様を説明したが、本発明に係るテキスト要約プログラムは、C D - R O M や D V D - R O M 等の記録媒体に記録されている形態で提供することも可能である。

30

【 0 0 6 0 】

本明細書に記載された全ての文献、特許出願及び技術規格は、個々の文献、特許出願及び技術規格が参照により取り込まれることが具体的かつ個々に記された場合と同程度に、本明細書中に参照により取り込まれる。

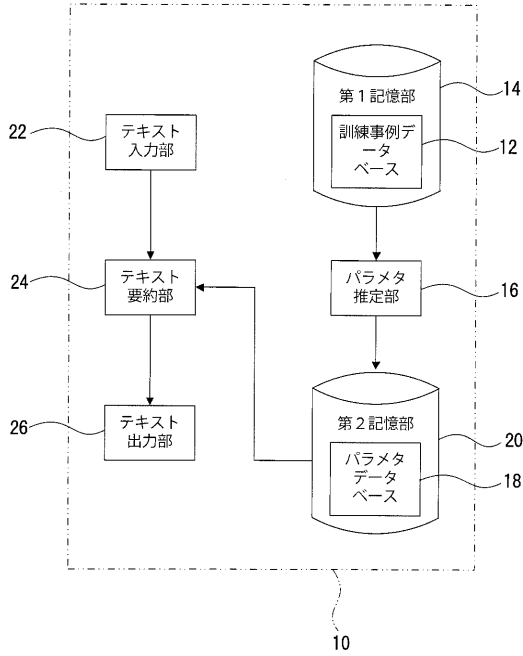
【符号の説明】

【 0 0 6 1 】

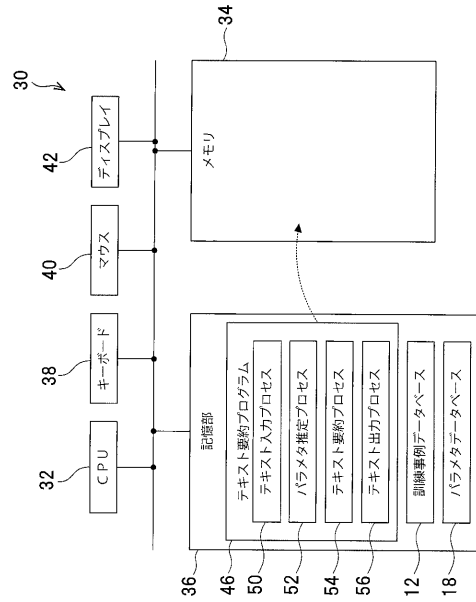
- 1 0 テキスト要約装置
- 1 2 訓練事例データベース
- 1 4 第 1 記憶部
- 1 6 パラメタ推定部
- 1 8 パラメタ・データベース
- 2 0 第 2 記憶部
- 2 2 テキスト入力部
- 2 4 テキスト要約部
- 2 6 テキスト出力部
- 3 0 コンピュータ
- 3 4 メモリ
- 3 6 記憶部
- 4 6 テキスト要約プログラム

40

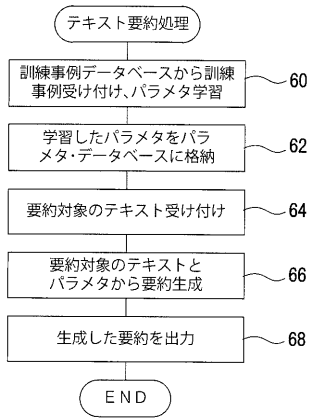
【 図 1 】



【 図 2 】



【 図 3 】



【 図 4 】

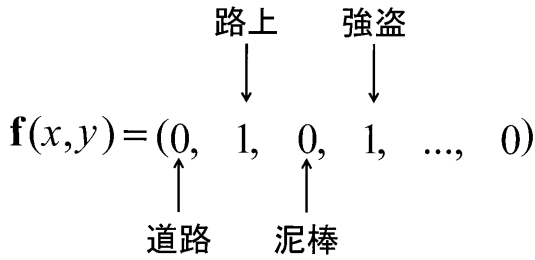
事例番号	テキスト	要約
1	9月26日午前、神奈川県横須賀市の路上で強盗事件が発生した。被害者は市内在住の70歳の男性。犯人は身長170センチメートル程度の中肉中背の男で、サングラスにマスク姿だったと言う。被害者の男性によれば男は原付バイクに乗って被害者に近づいたという。	9月26日午前、神奈川県横須賀市の路上で強盗事件が発生した。
2	9月27日未明、神奈川県横須賀市光の丘で火災が発生した。火元の民家は全焼した。火災当時居住者は外出しており建物内におらず、怪我人はいなかった。警察、消防は放火の可能性が高いとして調査を行っている。	9月27日未明、神奈川県横須賀市光の丘で火災が発生した。
...	...	...

【 図 5 】

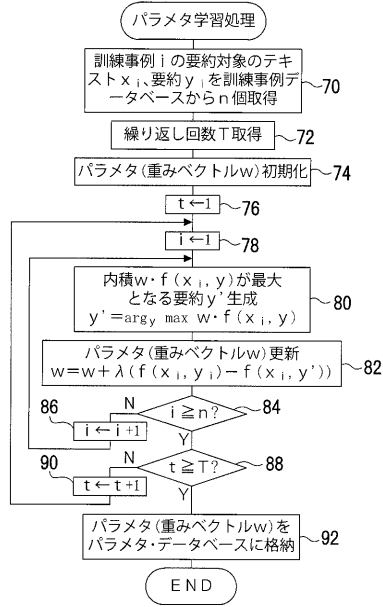
(A)

テキスト x	要約 y
9月26日午前、神奈川県横須賀市の路上で強盗事件が発生した。被害者は市内在住の70歳の男性。犯人は身長170センチメートル程度の中肉中背の男で、サングラスにマスク姿だったと言う。被害者の男性によれば男は原付バイクに乗って被害者に近づいたという。	9月26日に神奈川県横須賀市の路上で強盗事件が発生した。

(B)



【 図 6 】



【 図 7 】

重み	-0.03	0.51	-0.49	0.38	...
特徴	道路	路上	泥棒	強盗	...

w

【 図 8 】

表記	品詞	読み	標準形
9月	名詞:日時:連用	クガツ	9月
26日	名詞:日時:連用	ニジュウロクニチ	26日
午前	名詞:日時:連用	ゴゼン	午前
、	読点		、
神奈川県	名詞:固有:地	カナガワ	神奈川県
県	名詞:接尾辞:名詞	ケン	県
横須賀	名詞:固有:地	ヨコスカ	横須賀
市	名詞:接尾辞:名詞	シ	市
の	格助詞:連体	ノ	の
路上	名詞	ロジョウ	路上
で	格助詞:連用	デ	で
強盗	名詞	ゴウトウ	強盗
事件	名詞	ジケン	事件
が	格助詞:連用	ガ	が
発生	名詞:動作	ハツセイ	発生
し	動詞活用語尾	シ	し
...	...	...	...

【 図 9 】

