

(19) 日本国特許庁(JP)

(12) 特 許 公 報(B2)

(11) 特許番号

特許第3832613号

(P3832613)

(45) 発行日 平成18年10月11日(2006.10.11)

(24) 登録日 平成18年7月28日(2006.7.28)

(51) Int. Cl.

G06F 17/30 (2006.01)

F I

G06F 17/30 220A

G06F 17/30 170A

請求項の数 8 (全 14 頁)

(21) 出願番号	特願平10-270147	(73) 特許権者	000004352
(22) 出願日	平成10年9月24日(1998.9.24)		日本放送協会
(65) 公開番号	特開2000-99536(P2000-99536A)		東京都渋谷区神南2丁目2番1号
(43) 公開日	平成12年4月7日(2000.4.7)	(74) 代理人	100077481
審査請求日	平成16年4月8日(2004.4.8)		弁理士 谷 義一
特許権者において、実施許諾の用意がある。		(74) 代理人	100088915
			弁理士 阿部 和夫
		(74) 代理人	100105371
			弁理士 加古 進
		(72) 発明者	加藤 直人
			東京都世田谷区砧一丁目10番11号 日 本放送協会 放送技術研究所内
		審査官	辻本 泰隆
			最終頁に続く

(54) 【発明の名称】 自動要約装置および自動要約プログラムを記録した記録媒体

(57) 【特許請求の範囲】

【請求項1】

入力された原文を形態素解析し、原文に含まれる文字列を単語に分割しその品詞を付与した形態素情報を入力する形態素解析手段と、

入力された原文の要約率と入力された該原文の文字数とから削減すべき文字数の最小値を計算する文字削減数最小値計算手段と、

前記形態素解析手段から出力される前記形態素情報に基づいて、前記原文の単語列に対して先頭の単語列から順に、あらかじめ得られている置換知識との照合を行い、照合に成功した場合には置換する単語列を前記形態素情報に追加し、かつその置換コストを付与した単語ラティス構造を出力する置換知識検索手段と、

該置換知識検索手段から出力される前記単語ラティス構造を文末から文頭に向けて探索し、該単語ラティス構造の文末から各ノードまでの最大可能な文字削減数を計算する後向き文字削減数計算手段と、

前記置換知識検索手段から出力される前記単語ラティス構造を文頭から文末に向けて探索し、各パスに対して文頭からの文字削減数を計算し、その計算結果と前記後向き文字削減数計算手段から出力された文字削減数との和を求め、その和が所望の文字削減数よりも小さい場合には、そのパスを枝刈りし、その和が所望の文字削減数よりも小さくない場合には、そのパスを出力する前向き文字削減数計算手段と、

該前向き文字削減数計算手段から出力されるパスに対してその置換コストの和を計算し、文末に達したときには文頭から文末までのパスを出力する置換コスト計算手段と、

10

20

該置換コスト計算手段から出力されたパスの中で、文頭から文末までの置換コストの和が最小となるパスを求め、求めた最小パスの文字列を要約文として出力する置換コスト最小パス計算手段と

を有することを特徴とする自動要約装置。

【請求項 2】

前記置換知識検索手段は、パスごとに必須適用置換知識リストをもっておき、置換知識を適用する時にそのリストを参照することを特徴とする請求項 1 に記載の自動要約装置。

【請求項 3】

前記後向き文字削減数計算手段は、前記単語ラティス構造の文末ノード n (n は単語数) から文頭ノード 0 に向けてノード番号 i を 1 ずつ減少させて、各ノード i において、文末ノード n からその現ノード i までの最大可能な文字削減数すなわち文字削減数の最大値 $m_b(i)$ をダイナミックプログラミングにより求めることを特徴とする請求項 1 または 2 に記載の自動要約装置。

10

【請求項 4】

前記前向き文字削減数計算手段は、文末ノードから文頭ノードまでの文字削減数の最大値 $m_b(0)$ が、削減すべき所望の文字数の最小値 m よりも小さい場合には、入力された要約率では要約できない旨を出力することを特徴とする請求項 1 ないし 3 のいずれかに記載の自動要約装置。

【請求項 5】

前記前向き文字削減数計算手段は、文末ノードから文頭ノードまでの文字削減数の最大値 $m_b(0)$ が、削減すべき文字数の最小値 m よりも小さくない場合には、前記単語ラティス構造の文頭ノード 0 から文末ノード n に向けてノード番号 i を 1 ずつ増加させて、各ノード i までのすべてのパスにおいて、そのパスの文字削減数 $m_f(i)$ を計算し、前記後向き文字削減数計算手段で得られた前記ノード i における文字削減数 $m_b(i)$ との和をとり、その和 $m_f(i) + m_b(i)$ が削減すべき文字数の最小値 m よりも小さい場合には、そのパスを枝刈りし、その和 $m_f(i) + m_b(i)$ が削減すべき文字数の最小値 m よりも小さくない場合には、そのパスを出力することを特徴とする請求項 4 に記載の自動要約装置。

20

【請求項 6】

前記前向き文字削減数計算手段は、前向き計算の処理中に、正例がない単語に達したときに、条件(前向き文字削減数、必須適用置換リスト等)の同じ候補は、置換コストが最小でないパスも枝刈りすることを特徴とする請求項 5 に記載の自動要約装置。

30

【請求項 7】

コンピュータによって原文と所望の要約率から該要約率の中で最適な要約文を自動的に求めるための自動要約プログラムを記録した記録媒体であって、該自動要約プログラムはコンピュータに、

入力された原文を形態素解析させることで原文に含まれる文字列を単語に分割させ、かつその品詞を付与した形態素情報を生成させ、

入力された原文の要約率と入力された原文の文字数とから削減すべき文字数の最小値を計算させ、

40

前記形態素情報に基づいて、前記原文の単語列に対して先頭の単語列から順に、あらかじめ得られている置換知識との照合を行わせ、照合に成功した場合には置換する単語列を前記形態素情報に追加して、かつその置換コストを付与した単語ラティス構造を求めさせ、前記単語ラティス構造を文末から文頭に向けて探索させて、該単語ラティス構造の文末から各ノードまでの最大可能な文字削減数(後向き文字削減数)を計算させ、

前記単語ラティス構造を文頭から文末に向けて探索させて、各パスに対して文頭からの文字削減数を計算させ、その計算結果と前記後向き文字削減数との和を求め、その和が所望の文字削減数よりも小さい場合には、そのパスを枝刈りさせ、その和が所望の文字削減数よりも小さくない場合には、そのパスを選出させ、

該選出されたパスに対してその置換コストの和を計算させ、文末に達したときには文頭か

50

ら文末までのパスを出力させ、
該出力されたパスの中で、文頭から文末までの置換コストの和が最小となるパスを求め、
求めた最小パスの文字列を要約文として出力させることを特徴とする自動要約プログラム
を記録した記録媒体。

【請求項 8】

前記自動要約プログラムはコンピュータに、
文末ノードから文頭ノードまでの文字削減数の最大値が、削減すべき所望の文字数の最小
値よりも小さい場合には、入力された要約率では要約できない旨を出力させることを特徴
とする請求項 7 に記載の自動要約プログラムを記録した記録媒体。

【発明の詳細な説明】

10

【0001】

【発明の属する技術分野】

本発明は、原文に対して所望の要約率が与えられたときに、その所望の要約率の中で最適
な要約文を自動的に求める自動要約を行う自動要約装置および自動要約プログラムを記録
した記録媒体に関する。

【0002】

本発明は、TV ニュース等の文章を自動的に要約するのに好適であり、また字幕作成への
応用なども考えられる。

【0003】

【従来の技術】

20

自動要約とは、原文の単語列を短い単語列に置換することにより、原文を自動的に縮約す
ることである。自動要約を実現するためには、原文のどの単語列をどのような単語列に置
換するのかという置換知識が必要となる。例えば、次が置換知識の一例である。

【0004】

[置換知識 1] (原文単語列 要約文単語列)

[置換知識 1 a] 明らかにしました 表明, 置換コスト = 0.6

[置換知識 1 b] 明らかに 表明, 置換コスト = 0.3

[置換知識 1 c] まし (は空、省略を表す記号), 置換コスト = 0.2

ここで、置換知識に付随している「置換コスト」は、その置換知識を使う際のペナルテ
ィーを表しており、置換コストが 0 に近いほどその置換知識は使いやすいとしている。

30

【0005】

上述の置換知識 1 の例を使って、下記の原文 1 を次の要約率で要約することを考えてみよ
う。なお、要約率 = 要約文の文字数 / 原文の文字数 × 100 とする。

[原文 1] 「明らかにしました」 (8 文字)

[要約率 1] 要約率 60% 以下。

【0006】

要約率を 60% 以下にするとということは、削除する文字数でいうと次のようになる。

【0007】

8 文字 × (100 - 60)% = 3.2 文字

すなわち、原文中の 3.2 文字以上の文字数を削減しなければならない。

40

【0008】

従来の自動要約では、置換知識を出現順に順次適用することにより要約文を得ていた。ま
た、そこで使われる置換知識は人手で収集し作成していたために、その数は非常に少な
かった。(例えば、文献 1: 『山本ほか: 「文章内構造を複合的に利用した論説文要約シ
ステム GREEN」自然言語処理、Vol. 2, No. 2, pp. 39-55, 1994』
は、人手で作成した置換知識をはじめとする要約知識を順次適用することにより自動要
約している。)

【0009】

要約率 60% 以下の例では、上記文献 1 と同様に、削減文字数 3.2 文字以上を置換知識
1 の出現順に順次適用すると、置換知識 1 a が適用され、下記の要約文 1 a が得られる。

50

【 0 0 1 0 】

[要約文 1 a] 「表明」 (削減文字数 = 6 文字、置換コストの和 = 0 . 6)

【 0 0 1 1 】

【 発明が解決しようとする課題 】

しかし、最近、置換知識を自動的に作成する手法が開発され、大量の置換知識を簡単に得ることができるようになった。(例えば、文献 2 : 『加藤直人 : 「ニュース文要約のための局所的な要約知識獲得とその評価」電子情報通信学会言語理解とコミュニケーション研究会, N L C 9 8 - 1 6 , p p . 7 - 1 4 , 1 9 9 8 』)。

【 0 0 1 2 】

このようにして置換知識が大幅に増えると、今度は 1 つの原文の単語列に対して複数の置換知識を適用できる場合があり、そのため自動要約の際に、適用する置換知識間で競合が生じるようになって、最適な要約文が求められないという場合があった。

【 0 0 1 3 】

上述の例でも、置換知識 1 a を適用して得られる要約文 1 a の場合と、置換知識 1 b と 1 c を適用して得られる要約文 1 b の場合が競合するが、下記の要約文 1 b のほうが置換コストの和が 0 に近いので適切な要約となる。

【 0 0 1 4 】

[要約文 1 b] 「表明した」 (削減文字数 = 4 文字、置換コストの和 = 0 . 5)

そこで、本発明の目的は、上述のような点に鑑みて、原文と所望の要約率が与えられたときに、置換コストと文字削減数を利用して、置換知識を適切にかつ効率的に選択することにより、原文の最適な要約文を自動的に得ることを可能にすることにある。

【 0 0 1 5 】

【 課題を解決するための手段 】

上記目的を達成するため、請求項 1 の自動要約装置の発明は、入力された原文を形態素解析し、原文に含まれる文字列を単語に分割しその品詞を付与した形態素情報を入力する形態素解析手段と、入力された原文の要約率と入力された該原文の文字数とから削減すべき文字数の最小値を計算する文字削減数最小値計算手段と、前記形態素解析手段から出力される前記形態素情報に基づいて、前記原文の単語列に対して先頭の単語列から順に、あらかじめ得られている置換知識との照合を行い、照合に成功した場合には置換する単語列を前記形態素情報に追加し、かつその置換コストを付与した単語ラティス構造を出力する置換知識検索手段と、該置換知識検索手段から出力される前記単語ラティス構造を文末から文頭に向けて探索し、該単語ラティス構造の文末から各ノードまでの最大可能な文字削減数を計算する後向き文字削減数計算手段と、前記置換知識検索手段から出力される前記単語ラティス構造を文頭から文末に向けて探索し、各パスに対して文頭からの文字削減数を計算し、その計算結果と前記後向き文字削減数計算手段から出力された文字削減数との和を求め、その和が所望の文字削減数よりも小さい場合には、そのパスを枝刈りし、その和が所望の文字削減数よりも小さくない場合には、そのパスを出力する前向き文字削減数計算手段と、該前向き文字削減数計算手段から出力されるパスに対してその置換コストの和を計算し、文末に達したときには文頭から文末までのパスを出力する置換コスト計算手段と、該置換コスト計算手段から出力されたパスの中で、文頭から文末までの置換コストの和が最小となるパスを求め、求めた最小パスの文字列を要約文として出力する置換コスト最小パス計算手段とを有することを特徴とする。

【 0 0 1 6 】

ここで、好ましくは、前記置換知識検索手段は、パスごとに必須適用置換知識リストをもっておき、置換知識を適用する時にそのリストを参照する。

【 0 0 1 7 】

また、好ましくは、前記後向き文字削減数計算手段は、前記単語ラティス構造の文末ノード n (n は単語数) から文頭ノード 0 に向けてノード番号 i を 1 ずつ減少させて、各ノード i において、文末ノード n からその現ノード i までの最大可能な文字削減数すなわち文字削減数の最大値 m_i (i) をダイナミックプログラミングにより求める。

10

20

30

40

50

【0018】

また、好ましくは、前記前向き文字削減数計算手段は、文末ノードから文頭ノードまでの文字削減数の最大値 $m_0(0)$ が、削減すべき所望の文字数の最小値 m よりも小さい場合には、入力された要約率では要約できない旨を出力する。

【0019】

また、好ましくは、前記前向き文字削減数計算手段は、文末ノードから文頭ノードまでの文字削減数の最大値 $m_0(0)$ が、削減すべき文字数の最小値 m よりも小さくない場合には、前記単語ラティス構造の文頭ノード0から文末ノード n に向けてノード番号 i を1ずつ増加させて、各ノード i までのすべてのパスにおいて、そのパスの文字削減数 $m_f(i)$ を計算し、前記後向き文字削減数計算手段で得られた前記ノード i における文字削減数 $m_0(i)$ との和をとり、その和 $m_f(i) + m_0(i)$ が削減すべき文字数の最小値 m よりも小さい場合には、そのパスを枝刈りし、その和 $m_f(i) + m_0(i)$ が削減すべき文字数の最小値 m よりも小さくない場合には、そのパスを出力する。

10

【0020】

また、好ましくは、前記前向き文字削減数計算手段は、前向き計算の処理中に、正例がない単語に達したときに、条件(前向き文字削減数、必須適用置換リスト等)の同じ候補は、置換コストが最小でないパスも枝刈りする。

【0021】

上記目的を達成するため、請求項7の記録媒体の発明は、コンピュータによって原文と所望の要約率から該要約率の中で最適な要約文を自動的に求めるための自動要約プログラムを記録した記録媒体であって、該自動要約プログラムはコンピュータに、入力された原文を形態素解析させることで原文に含まれる文字列を単語に分割させ、かつその品詞を付与した形態素情報を生成させ、入力された原文の要約率と入力された原文の文字数とから削減すべき文字数の最小値を計算させ、前記形態素情報に基づいて、前記原文の単語列に対して先頭の単語列から順に、あらかじめ得られている置換知識との照合を行わせ、照合に成功した場合には置換する単語列を前記形態素情報に追加して、かつその置換コストを付与した単語ラティス構造を求めさせ、前記単語ラティス構造を文末から文頭に向けて探索させて、該単語ラティス構造の文末から各ノードまでの最大可能な文字削減数(後向き文字削減数)を計算させ、前記単語ラティス構造を文頭から文末に向けて探索させて、各パスに対して文頭からの文字削減数を計算させ、その計算結果と前記後向き文字削減数との和を求め、その和が所望の文字削減数よりも小さい場合には、そのパスを枝刈りさせ、その和が所望の文字削減数よりも小さくない場合には、そのパスを選出させ、該選出されたパスに対してその置換コストの和を計算させ、文末に達したときには文頭から文末までのパスを出力させ、該出力されたパスの中で、文頭から文末までの置換コストの和が最小となるパスを求め、求めた最小パスの文字列を要約文として出力させることを特徴とする。

20

30

【0022】

ここで、好ましくは、前記自動要約プログラムはコンピュータに、文末ノードから文頭ノードまでの文字削減数の最大値が、削減すべき所望の文字数の最小値よりも小さい場合には、入力された要約率では要約できない旨を出力させる。

【0023】

【発明の実施の形態】

本発明の実施の形態を説明するに先立ち、本発明に係る自動要約に必要な上述の置換コストについて説明する。

40

【0024】

自動要約に必要な要約知識は、置換知識と置換条件の2つから構成されている。置換知識は上述のように原文の単語列をどのような単語列に置換するかを規定する知識である。例えば、連体助詞の「の」を省略するという知識である。一方、置換条件とは置換知識の適用の良否を数値化したもの、すなわち上述の置換コストである。置換知識はその前後の単語列によって適用の良否が決まる。例えば、「日本の銀行」の「の/体助」を省略することはできない。

50

【0025】

そこで、置換コストは、置換知識の前後の単語列と、あらかじめ獲得しておいた置換条件との距離を計算している。すなわち、 i 番目から j 番目までの単語列 w_{ij} を、単語列 x_{ij} に置換するという置換コストを $\text{distsub}(w_{ij} \rightarrow x_{ij})$ と表すと、(1) 式で定義される(さらに詳しくは、上記文献2を参照。)

【0026】

【数1】

$$\text{distsub}(w_{ij} \rightarrow x_{ij}) = \begin{cases} g'(w_{ij} \rightarrow x_{ij}) & \text{if 正例があるとき} \\ 0.0 & \text{otherwise} \end{cases} \quad \dots (1)$$

ただし、

$$g'(w_{ij} \rightarrow x_{ij}) = (1.0 - g_{\text{low}}) \times g(w_{ij} \rightarrow x_{ij}, \text{正例}) + g_{\text{low}} \quad \dots (1a)$$

$$g_{\text{low}} = 0.01 \quad \dots (1b)$$

【0027】

上記(1)式は、正例がある場合には、

$g_{\text{low}} (= 0.01) \sim 1.0$ の値 ($0.0 \leq g(w_{ij} \rightarrow x_{ij}, \text{正例}) \leq 1.0$) を取り、0.0に近いほど置換することが可能であると定義されている。また、正例がない(適用される置換知識がない)場合には0.0を取る。

【0028】

また、本発明による自動要約アルゴリズムの概要を説明する。説明を簡単にするために、以下では1文を要約する場合を考える。複数の文にわたる場合には単純に連結すればよい。

【0029】

今、原文をある要約率以下に要約したいとする。このとき、 $m (= \text{原文の文字数} \times \text{要約率})$ 文字以上の文字を削除しなければならない。さらに、最適な要約であってほしい。ここで、「最適な要約」とは、適用した置換知識のコストの和(置換コスト)が最小となる場合であると定義する。したがって、自動要約とは、 m 文字以上の文字数を削除し、文頭から文末までの置換コストが最小のパス(最適パス)を求めることである。定式化すると、(2)式ようになる。

【0030】

【数2】

$$\text{argmin}_{x \in X} \sum_x \text{distsub}(w_{ij} \rightarrow x_{ij}) \quad (2)$$

$$X = \{ (x_0, \dots, x_{ij}, \dots, x_n) \mid \sum (|w_{ij}| - |x_{ij}|) \geq m \}$$

(2)式の解を求めるアルゴリズムについては図1、図2を用いて後述する。なお、本発明では、文字削減数と置換コストという2つの評価関数を用いているが、前者を計算する際にヒューリスティック関数(現在のノードからゴールまでの評価関数の予測値)を用いている。

【0031】

以下、図面を参照して本発明の実施形態を詳細に説明する。

【0032】

図1は、本発明の一実施形態の装置構成を示す。図1において、i1は原文を入力する端子であり、i2は要約率を入力する端子である。o1は要約文を出力する端子である。

【0033】

形態素解析装置1は、i1の端子に入力された原文を形態素解析し、原文に含まれる文字列を単語に分割しその品詞を付与した形態素情報を入力する。

【0034】

文字削減数最小値計算装置2は、i2の端子に入力された原文の要約率とi1の端子に入力された原文の文字数とから削減すべき文字数の最小値を計算する。置換知識検索装置3は、形態素解析装置1から出力される形態素情報に基づいて、原文の単語列に対して先頭の単語列から順に、あらかじめ得られている置換知識との照合を行い、照合に成功した場合は置換する単語列を上記形態素情報に追加し、かつその置換コストを付与した単語ラティス構造を出力する。本例では、その置換コストは、各置換知識に付随して置換知識と共に内部メモリ(図示しない)にあらかじめ格納されているものとする。

10

【0035】

後向き文字削減数計算装置4は、置換知識検索装置3から出力される単語ラティス構造を文末から文頭に向けて探索し、その単語ラティス構造の文末から各ノードまでの最大可能な文字削減数を計算する。

【0036】

前向き文字削減数計算装置5は、置換知識検索装置3から出力される単語ラティス構造を文頭から文末に向けて探索し、各パスに対して文頭からの文字削減数を計算し、その計算結果と後向き文字削減数計算装置4から出力された文字削減数との和を求め、その和が所望の文字削減数(入力された上記原文と要約率から算出)よりも小さい場合には、そのパスを枝刈りし、その和が所望の文字削減数よりも小さくない場合には、そのパスを出力する。

20

【0037】

置換コスト計算装置6は、前向き文字削減数計算装置5から出力されるパスに対してその置換コストの和を計算し、文末に達したときには文頭から文末までのパスを出力する。

【0038】

置換コスト最小パス計算装置7は、置換コスト計算装置6から出力されたパスの中で、文頭から文末までの置換コストの和が最小となるパスを求め、求めた最小パスの文字列を要約文としてo1の端子から出力する。

30

【0039】

図2は、図1の装置構成により自動要約の処理を行う手順の一例を示すフローチャートである。図2に従って、以下、本発明による自動要約の手順を説明する。

【0040】

まず、i1の端子に原文が入力されると、ステップS1では形態素解析装置1によって原文が形態素解析され、その単語分割と品詞が出力される。i2の端子に要約率が入力されると、次のステップS2では原文の文字数と要約率から削減すべき文字数の最小値が計算される。

40

【0041】

続くステップS3では、上記ステップS1で得られた形態素解析結果に基づいて、先頭の単語列から順に、あらかじめ得られている置換知識との照合を行い、照合に成功した場合には、形態素解析結果に置換する単語列を追加し、その置換コストを付与して単語ラティス構造を作成する。

【0042】

ステップS4~S7では、後向きの計算を行う。すなわち、ステップS4、S6、S7で単語ラティス構造の文末ノードn(nは単語数)から文頭ノード0に向けてノード番号iを1ずつ減少させて、各ノードiにおいてステップS5の処理を実行する。ステップS5では、文末ノードnからその現ノードiまでの最大可能な文字削減数(後ろ向き文字削減

50

数)すなわち文字削減数の最大値 $m_b(i)$ をダイナミックプログラミング(DP)により求め、保存する。

【0043】

次のステップS8では、文末ノードから文頭ノードまでの文字削減数の最大値 $m_b(0)$ が、削減すべき所望(指定の)の文字数の最小値 m よりも小さい場合には、ステップS9を実行した後、本自動要約処理を終了する。ステップS9では入力された要約率では要約できない旨を出力する。一方、文末ノードから文頭ノードまでの文字削減数の最大値 $m_b(0)$ が、削減すべき文字数の最小値 m よりも小さくない場合には、ステップS10を実行する。

【0044】

ステップS10～S16では前向きの計算を行う。ステップS10、S15、S16において単語ラティス構造の文頭ノード0から文末ノードnに向けてノード番号 i を1ずつ増加させて、各ノード i までのすべてのパスにおいてステップS11～S14の処理を実行する。まず、ステップS11では、そのパスの文字削減数(前向き文字削減数) $m_f(i)$ を計算し、上述のステップS5で求めたノード i における文字削減数 $m_b(i)$ との和を取り、ステップS12でその和 $m_f(i) + m_b(i)$ が削減すべき文字数の最小値 m よりも小さい場合には、このパスは最終的な解となり得ないので、ステップS13でそのパスを枝刈りする。文末にしたがい可能なパスの候補が増加していくが、このような枝刈りにより候補数を抑えることができる。一方、その和 $m_f(i) + m_b(i)$ が削減すべき文字数の最小値 m よりも小さくない場合には、ステップS14を実行する。ステップS14ではそのパスの置換コスト $cost(i)$ の和を求め、文末に達しているときには文頭から文末までのパスを出力する。

【0045】

次のステップS17では、上記のステップS14で求められたパスの中で、置換コストの和が最小となるパスを求め、要約文をo1の端子から出力する。

【0046】

さらに、図2の処理の具体的な一例を、下記の原文2を次のように要約する場合を用いて説明する。

【0047】

[原文2]

「福沢総理大臣は特別委員会で方針を明らかにしました」 (24文字)

[要約率2]

要約率70%以下。

【0048】

また、置換知識として、次があらかじめ得られているとする。

【0049】

10

20

30

[置換知識 2]

[置換知識 2 a] 総理大臣／普通名詞→首相／普通名詞,
置換コスト=0.01

[置換知識 2 b] 委員会／普通名詞→委／普通名詞,
置換コスト=0.02

[置換知識 2 c] 明らか／形容名詞 に／格助詞に し／サ変連用
まし／助動丁寧 た／助動過去
→表明／サ変名詞、置換コスト=0.6

[置換知識 2 d] 明らか／形容名詞 に／格助詞
→表明／サ変名詞、置換コスト=0.3

[置換知識 2 c] まし→／φ (φは空を表す記号)、置換コスト=0.2

ここで、単語は「表層表現／品詞」と表している。

【0050】

i 1の端子に原文が入力されると、ステップS 1では原文を形態素解析する。すると、原文 2はその形態素解析結果として、図3(a)に示すように、単語数(これをnと表す) 13個の単語に分割され、品詞が付与される。図3(a)において、上段の数字は単語間に文頭から順につけたノード番号である。このとき、文頭ノードの番号は0であり、文末ノードの番号は単語数n(=13)である。

【0051】

i 2の端子に要約率70%が入力されると、ステップS 2では、原文の文字数が24文字であることから、削減すべき文字数の最小値(これをmと表し、所望の文字削減数と呼ぶ)が次のように計算される。

【0052】

[所望の文字削減数]

$$m = 24 \text{文字} \times (100 - 70)\% = 7.2 \text{文字}$$

ステップS 3では、上記形態素解析結果に基づいて、先頭の単語列から順に置換知識2との照合を行い、照合に成功した場合には置換する単語列を形態素解析結果に追加し、その置換を付与した単語ラティス構造を作成する。すると、図3(b)に示すような単語ラティス構造が得られる。なお、図3(b)で、例えば、置換候補の「首相」の脇に記載した「0.01」は置換コストを表す。

【0053】

ステップS 4～S 7では後向きの計算を行う。ステップS 4, S 6, S 7で単語ラティス構造の文末ノード13から文頭ノード0に向けてノード番号iをi=13, 12, ..., 0と1ずつ減少させて、各ノードiにおいてステップS 4を実行する。ステップS 4では文末ノード13から現在着目しているノードiまでに最大可能な文字削減数(これを $m_b(i)$ と表し、後向き文字削減数と呼ぶ)を計算する。この際に、ダイナミックプログラミングにより効率的に処理を行う。すると、図3(c)に示すように、各ノードにおける後向き文字削減数が計算される。

【0054】

図3(c)において、例えば、ノード8での後向き文字削減数 $m_b(8)$ は、ケース8a 置換知識2d「まし」と置換知識2e「明らかに 表明」を適用。(文字削減数4)

ケース8b 置換知識2c「明らかにしました 表明」を適用。(文字削減数6)と2つの場合が考えられるが、ケース8bが文字削減数が最も大きいので、

10

20

30

40

50

$m_b(8) = 6$ と求められる。

【0055】

ステップS8では文末ノードから文頭ノードまでの文字削減数の最大値 $m_b(0) = 10$ が、所望の文字削減数 $m = 7.2$ も小さくない($m_b(0) = 10$ $m = 7.2$)ので、ステップS9は実行せずに、ステップS10を実行する。

【0056】

ステップS10～S16では前向き計算を行う。ステップS10, S15, S16で単語ラティス構造の文頭ノード0から文末ノード13に向けてノード番号 i を $i = 0, 1, \dots, 13$ と1ずつ増加させて、各ノード i におけるすべてのパスにおいてステップS11～S14を実行する。ステップS11では文頭から現在着目しているノード i までのそれぞれのパスにおける文字削減数(これを $m_f(i)$ と表し、前向き文字削減数と呼ぶ)を計算する。ノード5($i = 5$)の場合を例にとると、図4(a)に示すように、4つのパスが求められる。

10

【0057】

次に、ステップS5で求めた後向き文字削減数 $m_b(i)$ との和をとる。ノード5の例では、後向き文字削減数 $m_b(5) = 6$ との和をとると、図4(b)に示すようになる。

【0058】

図4(b)に示すその和($m_f(i) + m_b(i)$)が所望の文字削減数 m よりも小さい場合($m_f(i) + m_b(i) < m$)には、ステップS13を実行し、小さくない場合($m_f(i) + m_b(i) \geq m$)には、ステップS14を実行する。ステップS13ではそのパスを枝刈りする。ステップS14ではそのパスの置換コストの和を求める。ノード5の例では、パス5aは、前向き文字削減数と後向き文字削減数との和が所望の文字削減数($m = 7.2$)よりも小さいパスであるので、ステップS13で枝刈りされ、以降の処理では使われない。パス5b、5c、5dはその和が所望の文字削減数($m = 7.2$)よりも小さくないので、ステップS14が実行される。ステップS14では、それぞれのパスにおける置換コストの和を求める。ノード5の例では図4(c)に示すようになる。

20

【0059】

同様にして、ステップS14において文頭から文末までの各パスにおける置換コストの和が求められる。文末まで達すると、文頭から文末までのパスである、図5(a)に示すような、4つの要約候補が得られる。

30

【0060】

ステップS17ではその要約候補の中で、置換コストの和が最小であるパス13dが選択され、単語列をつないで得られる次の要約文がo1の端子から出力される(図5(b)参照)。

【0061】

[要約文2] 「福沢首相は特別委で方針を表明した」

(他の実施の形態)

以上の説明では適用される置換知識がそれぞれ独立であるとした。しかし、「総理大臣首相」のように、1度適用したら次回にも必ず適用しなければならない置換知識もある。これに対応するためには、パスごとに必須適用置換知識リストをもっておき、置換知識を適用する時にそのリストを参照する処理を図2のアルゴリズムに追加すればよい。

40

【0062】

また、上述したステップS10～S16の前向き計算の処理中には、正例がない単語(例えば、図4(a)の「で」)に達したときに、条件(前向き文字削減数、必須適用置換リスト等)の同じ候補は、置換コストが最小でないパス(例えば、図4(a)のパス5b)も枝刈りするという改善を加えることも可能である。

【0063】

なお、図1の形態素解析装置1、文字削減数最小値計算装置2等はROM等を利用したモジュール回路(デバイス)のようなハードウェアで構成しても、形態素解析ルーチン、文字削減数最小値計算ルーチンのようにコンピュータ制御で動作するソフトウェアで構成し

50

てもよい。また、本発明は、複数の機器（例えば、ホストコンピュータ、インターフェース機器、リーダー、プリンタなど）から構成されるシステムに適用しても、1つの機器からなる専用装置（例えば、自動要約装置、自動字幕作成装置など）に適用してもよい。

【0064】

また、本発明の目的は、前述した実施の形態の機能を実現するソフトウェアのプログラムコードを記録した記録媒体（記憶媒体）を、システムあるいは装置に供給し、そのシステムあるいは装置のコンピュータ（またはCPUやMPU）が記録媒体に格納されたプログラムコードを読み出し、実行することによっても、達成されることは言うまでもない。この場合、記録媒体から読み出されたプログラムコード自体が前述した実施の形態の機能を実現することになり、そのプログラムコードを記録した記録媒体（例えば、CD-ROM、MD、フロッピーなど）は本発明を構成することになる。

10

【0065】

【発明の効果】

以上の説明から明らかなように、本発明によれば、文字削減数と置換コストという2つの評価関数を用いて、与えられた要約率以下で、原文を最適に要約するので、原文と要約率を入力するだけで、最適な要約を自動的に求めることが可能となる。

【図面の簡単な説明】

【図1】本発明の一実施形態の自動要約装置の構成を示すブロック図である。

【図2】図1の装置の自動要約の処理手順を示すフローチャートである。

【図3】具体的な原文を入力した場合の本発明の一実施形態の各段階の処理内容と結果を順次に説明する図であり、(a)形態素解析結果、(b)単語ラティス構造、(c)後向き文字削減数計算を具体例で示す。

20

【図4】図3に連続する説明図であり、(a)前向き文字削減数の計算、(b)前向き文字削減数と後向き文字削減数との和の計算、(c)置換コストの和の計算を具体例で示す。

【図5】図4に連続する説明図であり、(a)要約候補、(b)要約文を具体例で示す。

【符号の説明】

- 1 形態素解析装置
- 2 文字削減数最小値計算装置
- 3 置換知識検索装置
- 4 後向き文字削減数計算装置
- 5 前向き文字削減数計算装置
- 6 置換コスト計算装置
- 7 置換コスト最小パス計算装置

30

【 図 5 】

(a) 【要約候補】

バス13a 福沢→総理大臣→は→特別→委→で→方針→を表明₂
(文字削減数=8, 置換コストの和=0.62)

バス13b 福沢→首相→は→特別→委員会→で→方針→を表明₂
(文字削減数=8, 置換コストの和=0.61)

バス13c 福沢→首相→は→特別→委→で→方針→を表明₂
(文字削減数=10, 置換コストの和=0.63)

バス13d 福沢→首相→は→特別→委→で→方針→を表明₁→した
(文字削減数=8, 置換コストの和=0.53)



(b) 【要約文】

出力：福沢首相は特別委で方針を表明した

フロントページの続き

(56)参考文献 特開平8 - 212228 (JP, A)

特開平10 - 63658 (JP, A)

加藤 直人, ニュース文要約のための局所的要約知識獲得とその評価, 情報処理学会研究報告, 日本, 社団法人 情報処理学会, 1998年 7月24日, 第98巻 第63号, 69~76

(58)調査した分野(Int.Cl., DB名)

G06F 17/30,

G06F 17/27