

クエリ依存文短縮と見出し生成への応用

西川 仁^{1,a)} 今村 賢治^{1,b)} 別所 克人^{1,c)} 牧野 俊朗^{1,d)} 松尾 義博^{1,e)}

概要: 文短縮は長い文を重要な部分のみを残して短く縮める処理であり、抽出型要約における重要な要素技術である。本稿では何らかのニュース記事に対して端的な見出しを付与する問題を考え、これをクエリ依存文短縮タスクとして定式化する。見出しの元とする文と、クエリの形で表現される短縮後も残されるべき情報の2つが与えられているものとし、後者をできるだけ保持したまま前者の係り受け木の枝刈り候補の中から良好な候補を探索する。これによって、より適切に見出しの元となる文の内容と文としての自然さを保った見出しが生成できることを示す。

1. はじめに

1つ以上のニュース記事からなるニュース記事集合に対して端的な要約文を付与する課題を考える。この課題は、典型的には、ヘッドライン生成（あるいは見出し生成）と呼ばれており、単一のニュース記事を対象とする場合 [2], [5], [20] とニュース・アグリゲーターなどによって収集された何らかのトピックに関する複数のニュース記事を対象とする場合 [1], [6] とがある。単一あるいは複数のニュース記事に対して適切な見出しを付与することができれば、ニュース・アグリゲーターを利用するユーザーの情報アクセスを補助することができる。そのため、我々と取り巻くテキストの量が増加するにしたがって、その重要性は増している。

ニュース記事集合に対する適切な見出しは、少なくとも2つの要件を満たす必要がある。1つは、見出しが本文の内容を適切に反映していることである。見出しは、通常、本文を読むべきか否かを判断する手がかりをユーザーに与える。そのため、本文中において重要とされる情報を見出しは保持していなければならない。もう1つは文として自然であることである。見出しも、当然、自然言語によって書かれた文の一種であることから、文として自然なものではなければならない。また、ユーザーに対して誤解を与えることのない、曖昧性の少ない文である必要もある。

見出しを生成するという見地に立った場合、多くの既存研究はこれらの要件を必ずしも十分に満たしているとはい

えない。まず、重要な情報が脱落する可能性がある。見出しを生成する際、多くのアプローチは、まず、大量の見出し候補を生成する。次に、それらの中から予め設定された目的関数を最大化するものを探索し、目的関数を最大化する候補を見出しとして出力する。その際、目的関数には、見出し候補の内容を評価する関数と言語としての自然さを評価する関数の2つが含まれており、これらが線形結合されたものが目的関数として用いられることが多い。このような場合、重要な情報を欠いている見出し候補であっても、言語としての自然さが著しく高く評価されている場合、最終的な見出しとして選ばれる可能性がある。また、もう1つの問題として、そもそも言語として自然な見出しを作ることも難しい。見出し候補の言語としての自然さを評価するには n-gram 言語モデルが使われることが多いが、統計的機械翻訳や音声認識の出力と同様に、n-gram 言語モデルのみに基づいた自然言語の生成は不可解な文を生成することがままある。

これらの問題に対し、本稿では、過去に提案されてきた手法群を踏まえ、いくつかのヒューリスティックを明示的に利用し、自然な見出しを生成することに主眼を置く。我々は以下に示すヒューリスティックを提案する。

- 見出しにおいて重要な役割を果たすと思われる名詞句（以降、クエリと呼ぶ）を事前に同定しておき、見出しを生成する際にクエリが見出しにできる限り保存されるように陽に制約を与える。これによって生成された見出しが重要な情報を欠いているという危険を回避することができる。
- 一からの自然言語生成は行わず、本文から、本文を構成する文を1つ選び出し、その文を短縮することで見出しとする。これによって自然言語生成器が不可解な

¹ 日本電信電話株式会社 NTT メディアインテリジェンス研究所

a) nishikawa.hitoshi@lab.ntt.co.jp

b) imamura.kenji@lab.ntt.co.jp

c) bessho.katsuji@lab.ntt.co.jp

d) makino.toshiro@lab.ntt.co.jp

e) matsuo.yoshihiro@lab.ntt.co.jp

原文	東海道新幹線で、6年ぶりとなる新型車両「N700A」が、8日から営業運転を開始し、東京駅と新大阪駅で出発式が行われました。
クエリ	新幹線、新型車両、出発式
見出し	東海道新幹線で、新型車両「N700A」が、営業運転を開始し、新大阪駅で出発式が行われました。

図1 入力となる文とクエリ、および本稿で述べる手法によって生成された見出しの例。

文を生成する危険を回避することができる。

- 文を短縮する際には、文を係り受け解析することによって得られる係り受け木の枝刈りを行うことで文を短縮する。これによって、係り受け解析が正しくなされている限りは、文法的でない見出しが生成される危険を回避することができる。
- 係り受け木の枝刈りを行う際には、述語とその必須格の組が保存されるように枝刈りを行う。また、生成された見出しにおいても約物が正しく用いられるよう制約を加える。これによって不自然な見出しが生成される危険を回避することができる。

本稿で扱う課題は、見出しの元となる文とクエリが与えられている条件の下で、上に述べた制約を満たす部分木のうち最も見出しとして良好なものを探索する課題として定式化される。これはクエリ依存文短縮課題として表現することができ、本稿では見出し生成をクエリ依存文短縮課題として扱う。本稿は、上に述べた工夫により、より内容に富み、かつ自然な見出しが生成できることを示す。

本稿は以下のように構成されている。次の2節では関連研究を整理し、本稿のアプローチの妥当性を示す。3節では提案手法について詳しく述べる。4節ではモデルが必要とするパラメータを得る方法について述べる。5節では実際に見出しを生成する処理であるデコードの方法について述べる。6節では提案する手法を評価する実験について述べる。7節では実験の結果について述べ、議論を行う。8節では本稿をまとめると共に、今後の課題について議論する。

2. 関連研究

文書の見出しを生成する方法はいくつも提案されているが、大まかに以下の4つの種類に分類できる。

2.1 部分文字列の抽出

元の文書を構成する文字列のなかから、妥当な部分を抽出することによって見出しとすることができる。Filippova [6]らは複数の文書が与えられている状況において、それらの文書集合に対して端的な見出しを生成する課題を扱っている。Filippovaらは入力となる複数の単語列から単語をノードとする有向グラフを作り、この中から多くの単語

列で出現する経路を探索することで端的な見出しを生成する方法を提案している。見出しを生成すること目的として作られたものではないが、平尾ら [22]は、文の短縮を系列ラベリングとして捉え、文を構成する単語のうち、見出しに残すべきものを同定することで文を短縮している。見出しの元として適切な文を選び出すことができれば、平尾らによる方法で端的な要約文を生成することができる。

これらの手法の強みは統語的な解析を必要としないため、全体的な処理が軽量になること、また統語的な解析の誤りに対して頑健であることである。一方で、文の自然性の担保を n-gram に基づく言語モデルに強く依存しているため、文の自然性の維持が難しい。

2.2 構文木の枝刈り

統語的な情報によって得られる構文木を枝刈りすることによっても文を短く縮めることができる。枝刈り候補の中から最良の候補を選択する方法については、規則に基づくもの [9], [15], [18], n-gram 言語モデルなどの統計に基づくもの [3], [10], [23], 述語項構造に基づくもの [12], [17], 同期文法に基づくもの [4], [16]がある。

これらの手法は正しい統語解析の結果が与えられる限りは文法性が著しく損なわれることがなく、文としての自然さを保つという点において頑健である。一方、事前に統語的な解析が必要であるため、弱みは計算量である。

2.3 パタンとスロットの穴埋め

事前にスロットを持つパタンを数多く用意しておき、新しい文書が与えられた際に適切なパタンを選び出し、文書中からスロットを埋める要素を抽出してパタンにはめ込むことで見出しを生成する方法もある。Alfonseca ら [1]は大規模なニュース記事集合から見出しの原型となるパタンをあらかじめ抽出しておき、新しい文書が与えられた際にパタンのスロットを埋めることで見出しを生成している。

2.1節および2.2節で述べた方法では与えられた文書に含まれない表現を利用した見出しを生成できないのに対し、この方法は柔軟に様々な見出しを生成できるというメリットがある。その一方、事前にパタンを大量に用意しなければならないこと、事前に用意されているパタン以外で見出しが生成できないこと、またパタンのスロットを埋めるための情報抽出を高精度で行う必要があるといったデメリットがある。

2.4 統計的機械翻訳

最後の方法は、統計的機械翻訳 [11]に似たもので、まず、入力文書を単語の集合に分解する。次に、それらの単語の中から、単語が持つ重要度の和と単語の順序に与えられる言語尤度が最も高くなる単語の順列を選び出す。

Banko [2]らは単語の重みを tf-idf [14]で、言語尤度を

n-gram 言語モデルでそれぞれ与え、ビームサーチで探索を行っている。廣嶋 [20] らは言語モデルを用いるところは変わらないが、単語の重みづけに Support Vector Machine を利用している。Deshpande [5] は tf-idf と n-gram 言語モデルを用いる点は変わらないが、デコードに工夫しており、乱択法を用いることでビームサーチに比べ高速に最適解を発見する方法を提案している。

この方法の問題点は2つある。1つは文の自然性の担保の難しさである。この方法は、文の自然性の担保を n-gram 言語モデルに完全に依存している。そのため、統計的機械翻訳や音声認識の出力同様、不自然な文を出力することがままあり、品質の担保が難しい。もう1つはデコードの難しさである。この、単語を順列のなかからそれらの重要度と言語尤度に基づいて最良のものを探索する問題は NP 困難であり [5], 良好な解を素早く探し出すのは容易ではない。

このように、それぞれのアプローチに一長一短がある。本稿では、2.2 節で述べた、構文木の枝刈りに基づく方法を取る。本稿で扱う入力には主にニュース記事であることから、係り受けで一定の精度が得られることが期待できる。そのため、係り受け解析の結果を利用することで自然な文を生成できることが期待できる。また、2.1 節で述べた方法ほどではないが、2.3 節や 2.4 節で述べた方法と比べ比較的処理が軽量であるというのメリットである。

3. モデル

入力された、 n 個の文節からなる係り受け木を t とする。クエリとして入力された m 個の名詞句の集合を $Q = \{q_0, q_1, \dots, q_{m-1}\}$ とする。係り受け木 t の部分木をなすもののうち、後述するヒューリスティックを満たす部分木の集合を S とし、その要素を $s \in S$ とする。部分木 s に含まれるクエリの数を返す関数を $g: s, Q \rightarrow \mathbb{N}_0^+$ とする。 \mathbb{N}_0^+ は 0 を含む自然数である。また、部分木 s に対して、事前に定めたパラメータ W に基づき、言語としての自然さを与える関数を $f: s, W \rightarrow \mathbb{R}$ とする。 \mathbb{R} は実数である。部分木 s に含まれる文節の長さの合計を返す関数を $l: s \rightarrow \mathbb{N}_0^+$ とする。 k を、見出しの最大長とする。

3.1 目的関数

このとき、目的関数を以下のように定義する。

$$\hat{s} = \arg \max_{s \in S'} f(s, W) \quad (1)$$

$$S' = \arg \max_{s \in S} g(s, Q) \quad (2)$$

s.t. $l(s) \leq k$

ここで、 S' は、ヒューリスティックを満たす部分木の集合 S のうち、クエリの集合 Q に含まれるクエリを含む

数が最も多い、 S の部分集合 $S' \subseteq S$ であり、 \hat{s} は S' に含まれる要素のうち関数 f が最も高い値を与えるものである。すなわち、 k によって与えられる長さの制約とヒューリスティックを満たす部分木の集合 S のうち、クエリの集合を出来る限り多く含むものの中で、最も自然なものを見出しとして採用する。

3.2 ヒューリスティック

ヒューリスティックを満たす部分木 $s \in S$ の自然性を保つため、種々のヒューリスティックを加える。

3.2.1 根

部分木 $s \in S$ の根は、以下のいずれかの条件を満たすものとする。

- 元の木 t の最後の文節であること。
- s の最後の文節（根）の内容部の主辞が動詞であること。
- s の最後の文節（根）の内容部の主辞がサ変名詞であること。

3.2.2 必須格

必須格を伴わない述語の存在は文の自然性を著しく損なう。そのため、述語が見出しに含まれる際にはその必須格も見出しに含める必要がある。

必須格の同定には述語とその必須格を格納した辞書を用意し、これを利用した。辞書は、以下の手続きで作成された：

- (1) 約 23 億文を収集し、これを係り受け解析する。
- (2) 係り受け解析結果から述語とその格要素を集計する。集計の際には以下の処理を行う。
 - (a) 係り受け解析の結果から述語を取り出す。このとき、受身は除外する。
 - (b) 述語に直接係る文節のうち、機能部が格助詞である要素を取り出す。
 - (c) 1 つ以上の有効格を持つものを、述語毎に集計する。
- (3) 述語と格情報をフィルタリングする。以下の2つの条件を満たしているもののみを辞書に掲載する。
 - 50 回以上出現する述語であること。
 - 述語と X 格のペアが、対数尤度比検定において、危険率 0.1% で有意に多く共起しており、かつ、述語全体における、平均の X 格出現率よりその述語の X 格出現率が 10% 以上高いこと。

すなわち、述語とその格要素の組のうち、コーパス中に頻出し、かつ平均より有意に多く出現する組を辞書に格納したことになる。

3.2.3 形式名詞

「こと」「の」などの形式名詞は自立語として独立した1つの文節をなすが、一方で形式名詞が導く節がなければ文は意味をなさない。そこで、形式名詞が見出しに含まれる

場合、それに係る文節も見出しに含まれるようにする。

3.2.4 括弧

入力される文には括弧が含まれていることがままあり、括弧の片方のみが含まれる見出しは文の自然性を著しく損なう。そこで、後述するデコードにおいて、最後に見出しを出力する際、開括弧と閉括弧が含まれている数を検査し、これが等しくない場合はその見出し候補は出力せず、次に目的関数値が高い見出し候補の括弧の数を検査し、開括弧と閉括弧の数が等しければそちらを出力する。

4. パラメータ

ここでは、文としての自然さを評価する関数 f について述べる。関数 f は n -gram に基づく言語尤度と、係り受けに基づく言語尤度によって、見出しの言語尤度を計算するものとし、そのパラメータは以下のように求める。

4.1 n -gram に基づく言語尤度

n -gram による言語尤度は文の自然性を評価する最も基本的な方法である。本稿でもこれを利用し、見出し候補の中から文として自然なものを見つけ出すために用いる。1991年から2007年までの17年分の毎日新聞記事データ集^{*1}から3-gram言語モデルを構築し、これによって見出し候補の言語尤度を評価する。

4.2 係り受けに基づく言語尤度

ある文節がある文節に係る確率を計算することでも文の尤度を計算することができる。係り受け言語モデルでは、係り元の文節の内容部の主辞と機能部の主辞、係り先の文節の主辞を3-gramとみなして言語モデルを構築する。例えば、「計画を発表した」という文があったとする。この文は「計画を」と「発表する」の2つの文節からなっており、この文の言語尤度は $p(\text{発表} | \text{計画}, \text{を})$ という3-gramで求められる。構築の際には、4.1節と同様に、1991年から2007年までの17年分の毎日新聞記事データ集から言語モデルを構築した。

5. デコード

本稿では、ビームサーチを用いて探索を行う。これは、 n -gramによる言語尤度を利用することで探索空間が著しく複雑になり、課題を整数計画問題として表現しソルバーを用いて解を求めることが困難であるためである。

上に示した目的関数は、式(1)の中に式(2)が入れ子になっており、最適化が難しい。そこで、目的関数を以下のように書き換える。

$$\hat{s} = \arg \max_{s \in S} \{f(s, W) + \lambda g(s, Q)\} \quad (3)$$
$$\text{s.t. } l(s) \leq k$$

式(3)は、式(2)が式(1)に組み込まれた形になっていることに注意されたい。 λ はクエリが見出しに含まれていることに対する重みである。これに大きな値を与えておくことで、式(1)と式(2)の最適化を逐次的に行った場合と同様の結果が得られる。

探索にあたっては、 k で与えられる長さに関する制約と上述のヒューリスティックとを満たす部分木をその目的関数値と共に仮説 h として保持しておく。通常、 h は複数存在し、その集合を H とする。初期仮説集合 H^0 から探索を始め、仮説集合に含まれる仮説を逐次的に拡張するが、その過程で、目的関数値の低い仮説については探索を打ち切る。最後に得られた仮説集合のうち、目的関数値が最大の仮説を見出しと出力する。

デコードの大きな手続きは以下の通りである：

- (1) 根の候補となる文節を列挙しておき、それらを初期仮説集合 H^0 とする^{*2}。なお、文の途中にある述語を根とする仮説に対しても正しい言語尤度を計算できるように、それらを初期仮説集合に加える際は、述語を終止形にした上で言語尤度を計算するようにする。
- (2) H^0 に含まれる仮説 $h \in H^0$ を1つずつ取り上げ、それに係る文節を h に加える。その際、長さに関する制約およびヒューリスティックを満たすかどうか検査し、満たす場合は目的関数値を更新した上で次の仮説集合 H^1 に加える。
- (3) 仮説集合 H^i から次の仮説集合 H^{i+1} を作成するとき、 H^i に含まれる仮説のうち上位 b 個のみを取り上げる(ビームサーチ)。
- (4) H^{n-1} までの仮説集合を作成し終わったのち、もっとも目的関数値が高い仮説 $h^* \in \{H^0, H^1, \dots, H^{n-1}\}$ を見出しとして出力する。

実装にあたっては、 n 個の、サイズ b の優先度つきキューを用意し、優先度として目的関数値を与えて仮説を保持しておけばよい。仮説を展開する際には、仮説が上述の制約を満たしているか検査し、制約を満たしていれば目的関数値を更新した上で次のキューに移し替える。これを n 回繰り返すことで、上記の手続きは自然な形で実装することができる。

6. 実験

本節では、提案した手法を評価するための実験について述べる。実験では、形態素解析には測らによる形態素解析器[7]を、係り受け解析には今村らによる係り受け解析

^{*1} <http://www.nichigai.co.jp/sales/corpus.html>

^{*2} すなわち、ここで述べる探索においては、与えられている木の枝を刈るのではなく、木の根となりうる文節から枝を段階的に伸ばしていくことになる。

器 [8] を利用した。デコーダは Perl で独自に実装した。

6.1 データ

評価のため、ウェブ上から 250 種類のニュース記事を集めた。250 記事それぞれから、見出しの元としてふさわしいと思われる文を 1 つずつ人手で抜き出した。別途、当該記事において重要と思われる名詞句を、1 つの記事に対して最大 3 つまで抜き出した。抜き出された名詞句はクエリとして利用される。それぞれのニュース記事に対して 4 人の作業員が見出しを作成した。作業員には、ニュース記事から抽出された文に新しい表現を加えることはせず、文に含まれる不要と思われる部分を削除することで見出しを作成するように指示を与えた。見出しの長さは 10 文字から 30 文字程度とし、最大でも 50 文字以内との指示を与えた。人手によって作成された見出しは後述する BLEU を用いた評価に用いる。

6.2 ハイパーパラメータ

デコードの際には 3 種類のパラメータ λ と b 、 k を設定する必要がある。 λ はクエリが見出しに含まれている際に目的関数に与えられる重みである。 λ の値は 4.1 節および 4.2 節で説明した n-gram 言語モデルおよび係り受け言語モデルの対数尤度^{*3} に対して十分大きければよく、100 とした。ビーム幅 b は、解の精度と計算に要する時間とを考慮し、100 とした。見出しとなる文の長さの上限 k は 50 文字とした。

6.3 評価

冒頭で述べたように、本稿の目的は、何らかの文書に対して端的な要約文を付与することである。このことから、少なくとも以下の点を評価する必要がある。

- 見出しの内容性。見出しが、元となる文書の内容を適切に反映しているかどうか。
- 見出しの文としての自然さ。見出しが、文法的で、読み手に対して違和感を与えないものになっているかどうか。

これらを実験するため、BLEU [13] による評価と人手による評価の 2 つを行う。見出しの内容性については、BLEU を用いて、人手によって作成された見出しと比較することで評価する。見出しの文としての自然さについては、被験者に見出しを評価させることで評価する。

6.3.1 BLEU

要約の内容性を測るため、BLEU [13] を用いる。BLEU は 1 つの文を対象とする評価法であることに加え、複数の参照訳 (参照要約) がある状況を想定して設計されており、複数の参照要約がある今回のデータに対して自然に適用で

表 1 評価者への指示。

得点	基準
5 点	文に違和感をまったく感じない。問題なく文意をくみ取ることができる。
4 点	多少違和感を感じる文であるが、文意をくみ取ることができる。
3 点	違和感のある文ではあるが、なんとか文意をくみ取ることができる。
2 点	かなり不自然な文であり、文意をくみ取ること困難さを感じる。
1 点	非常に不自然な文であり、文意をくみ取ることができない (意味不明)。

きる。BLEU は平尾らによる日本語を対象とした文短縮の論文 [21], [22] でも評価に利用されている。

評価用のプログラムは Papineni らによる文献 [13] に従って独自に実装した。

6.3.2 人手による評価

要約の自然性を測るため、1 人の日本語母語話者に要約文を提示し、要約文を評価させた。被験者に見出しを与え、1 から 5 までの得点を与えるよう指示した。評価の基準は表 1 に示す。

6.4 ベースライン

実験では以下の 3 種類の手法を比較する。

- 統計 + クエリ制約 + ヒューリスティック (手法 1) 提案手法。n-gram 言語モデルおよび係り受け言語モデルで自然性を維持しつつ、長さに関する制約を満たす、クエリができる限り含まれる見出しを探索する。その際、前述のヒューリスティックも加味する。これを手法 1 と呼ぶ。
- 統計 + クエリ制約 (手法 2) 提案手法。n-gram 言語モデルおよび係り受け言語モデルで自然性を維持しつつ、長さに関する制約を満たす、クエリができる限り含まれる見出しを探索する。ただし、ヒューリスティックは加味しない。これを手法 2 と呼ぶ。
- 統計 (手法 3) ベースライン。n-gram 言語モデルおよび係り受け言語モデルで自然性を維持しつつ、長さに関する制約を満たす見出しを探索する。見出しの内容性については、制約の形ではなく、見出しが含む内容語^{*4} の tf-idf 値で与える。これを手法 3 と呼ぶ。tf-idf 値については 4.1 節で述べた n-gram 言語モデルに基づく言語尤度と同様、1991 年から 2007 年までの 17 年分の毎日新聞記事データ集から計算した。

これらの比較によって明らかにしたいことは以下の通りである。

- クエリ制約の効果。クエリを陽に制約として加え、できる限り見出しに残すようにすることで、どの程度内

*3 通常、負の値を取る。

*4 名詞、動詞、形容詞および未知語。

表 2 BLEU による評価の結果.

	BLEU
手法 1	0.492
手法 2	0.380
手法 3	0.313

表 3 人手による評価の結果.

	スコア
手法 1	3.57
手法 2	2.98
手法 3	2.22

容性が改善されるか.

- ヒューリスティックの効果. ヒューリスティックによって, 見出しの文としての自然さがどの程度改善されるか.

7. 結果と考察

7.1 BLEU

BLEU を用いて評価した結果を表 2 に示す. 手法 1, 手法 2, 手法 3 の順に良好な結果を得た. これらの手法の間にはすべて有意な差が認められた.

まず, 手法 2 と手法 3 を比べると, クエリを陽に制約として与えることによって BLEU が大幅に改善されることがわかる. さらに, 手法 1 と手法 2 を比べると, ヒューリスティックを与えることでも BLEU が向上することがわかる. 人手で作成された見出しを見ると, これらにおいては述語項構造が適切に見出しの中に保存されている. そのため, ヒューリスティックを用いることで人手で作成した見出しと同様に述語項構造が見出しの中に保存されることになり, これがヒューリスティックによって BLEU が改善された理由であると考えられる.

手法 1 においても残存する誤りみると, 係り受け解析の誤りによるものがその多くを占めた. 本稿において提案している, 係り受け木を枝刈りすることで見出しを生成する方法では, 係り受け解析の誤りに非常に敏感であるという問題点がある.

7.2 人手による評価

人手による評価の結果を表 3 に示す. BLEU による評価と同様, 手法 1, 手法 2, 手法 3 の順に良好な結果を得た.

まず, 手法 3 によって生成された見出しの人手による評価の結果を分析すると, 必須格を欠いている述語を含む文が多く見られる. このような文は正しく文意をくみ取ることができず, 人手による評価においては大きなデメリットとなる. また, 入力された文に含まれる括弧が片方だけ残っている場合があるなど, 文としての自然さを欠くものが多く存在しており, これらが人手による評価において手法 3 が劣後した理由と考えられる.

手法 2 は, 手法 3 にクエリ制約が加えられたもので, 文の自然さに直接寄与する改良が加えられたものではないが, 手法 2 に比べ評価が改善している. 手法 2 によって生成された見出しは当然クエリとなっている名詞句が文に多く残っている. これらは述語の主格や目的格となっていることが多く, 結果として, 必須格を欠いている述語の数が減少している. このことが手法 2 が手法 3 よりよい結果を得ている主たる理由と考えられる.

最後に手法 1 のよって生成された見出しのうち, 低い評価を得たものを見ると, 第一の理由は必須格辞書の網羅性が挙げられる. いくつかの動詞が辞書に含まれておらず, そのためそれらの動詞が必須格を欠いた状態で見出しに含まれており, そのため不自然な文が生成されている. 特に, いわゆる事態性名詞と呼ばれる名詞を述語とする述語項構造 [19] については, 現在の辞書に基づく方法では対応することができない. 事態性名詞を述語として持つ述語項構造の必須格が脱落することによって文意が失われている場合が見られる. こういった問題に対しては辞書ではなく陽に述語項構造解析器を用いることが必要である.

8. まとめと今後の課題

本稿では, 単一のニュース記事に含まれる文と, 重要な情報を表現したクエリを入力として, 入力された文の係り受け木をクエリを残すように刈り込むことで簡潔な文を生成する課題を扱った.

既存の手法群と異なり, クエリを陽に制約として扱い, できる限りクエリが見出しから脱落しないようにすること, またいくつかのヒューリスティックを加えて係り受け木の刈り込みに制限をかけることで, 内容性, 自然性いずれの観点においても良好な見出しを生成することができた.

残された課題として, まず係り受け解析の結果により頑健な方法を考える必要がある. これは解析木から n-best 解を受け取ることで対処できる可能性がある.

また, 事態性名詞を述部を持つ述語項構造を見出しで維持するためには述語項構造解析の利用が不可欠であり, こちらについても利用する必要がある.

また, より文を短くするためには何らかの言い換え処理が不可欠となるため, これらの獲得および処理についても扱う予定である.

参考文献

- [1] Alfonseca, E., Pighin, D. and Garrido, G.: HEADY: News headline abstraction through event pattern clustering, *Proceedings of the 51st Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pp. 1243–1253 (2013).
- [2] Banko, M., Mittal, V. O. and Witbrock, M. J.: Headline generation based on statistical translation, *ACL '00: Proceedings of the 38th Annual Meeting on Association for Computational Linguistics*, pp. 318–325 (2000).

- [3] Clarke, J. and Lapata, M.: Constraint-based Sentence Compression An Integer Programming Approach, *Proceedings of the 21st International Conference on Computational Linguistics and 44th Annual Meeting of the Association for Computational Linguistics (COLING-ACL)*, pp. 144–151 (2006).
- [4] Cohn, T. and Lapata, M.: Sentence Compression as Tree Transduction, *Journal of Artificial Intelligence Research*, Vol. 34, pp. 637–674 (2009).
- [5] Deshpande, P., Barzilay, R. and Karger, D.: Randomized Decoding for Selection-and-Ordering Problems, *Human Language Technologies 2007: The Conference of the North American Chapter of the Association for Computational Linguistics; Proceedings of the Main Conference*, pp. 444–451 (2007).
- [6] Filippova, K.: Multi-Sentence Compression: Finding Shortest Paths in Word Graphs, *Proceedings of the 23rd International Conference on Computational Linguistics (Coling 2010)*, pp. 322–330 (2010).
- [7] Fuchi, T. and Takagi, S.: Japanese morphological analyzer using word co-occurrence: JTAG, *Proceedings of the 36th Annual Meeting of the Association for Computational Linguistics and the 17th International Conference on Computational Linguistics (ACL-COLING)*, pp. 409–413 (1998).
- [8] Imamura, K., Kikui, G. and Yasuda, N.: Japanese Dependency Parsing Using Sequential Labeling for Semi-spoken Language, *Proceedings of the 45th Annual Meeting of the Association for Computational Linguistics Companion Volume Proceedings of the Demo and Poster Sessions*, pp. 225–228 (2007).
- [9] Jing, H.: Sentence Reduction for Automatic Text Summarization, *Proceedings of the Sixth Conference on Applied Natural Language Processing (ANLP)*, pp. 310–315 (2000).
- [10] Knight, K. and Marcu, D.: Summarization beyond sentence extraction: A probabilistic approach to sentence compression, *Artificial Intelligence*, Vol. 1, No. 139, pp. 91–107 (2002).
- [11] Koehn, P.: *Statistical Machine Translation*, Cambridge University Press (2010).
- [12] Morita, H., Sasano, R., Takamura, H. and Okumura, M.: Subtree Extractive Summarization via Submodular Maximization, *Proceedings of the 51st Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pp. 1023–1032 (2013).
- [13] Papineni, K., Roukos, S., Ward, T. and Zhu, W.-J.: Bleu: a Method for Automatic Evaluation of Machine Translation, *Proceedings of 40th Annual Meeting of the Association for Computational Linguistics*, pp. 311–318 (2002).
- [14] Sparck Jones, K.: A statistical interpretation of term specificity and its application in retrieval, *Journal of Documentation*, Vol. 28, pp. 11–21 (1972).
- [15] Toutanova, K., Brockett, C., Gamon, M., Jagarlamudi, J., Suzuki, H. and Vanderwende, L.: The PYPHY Summarization System: Microsoft Research at DUC2007, *Proceedings of The Document Understanding Conference (2007)*.
- [16] Woodsend, K., Feng, Y. and Lapata, M.: Title Generation with Quasi-Synchronous Grammar, *Proceedings of the 2010 Conference on Empirical Methods in Natural Language Processing (EMNLP)*, pp. 513–523 (2010).
- [17] Yoshikawa, K., Iida, R., Hirao, T. and Okumura, M.: Sentence Compression with Semantic Role Constraints, *Proceedings of the 50th Annual Meeting of the Association for Computational Linguistics (Volume 2: Short Papers)*, pp. 349–353 (2012).
- [18] Zajic, D. M., Dorr, B. J., Lin, J. and Richard, S.: Multi-Candidate Reduction: Sentence Compression as a Tool for Document Summarization Tasks, *Information Processing and Management*, Vol. 43, pp. 1549–1570 (2007).
- [19] 小町 守, 飯田 龍, 乾健太郎, 松本裕治: 名詞句の語彙統語パターンを用いた事態性名詞の項構造解析, *自然言語処理*, Vol. 17, No. 1, pp. 141–159 (2010).
- [20] 廣嶋伸章, 長谷川隆明, 奥 雅博: Web ページのヘッドライン生成のための統計的要約, *自然言語処理*, Vol. 12, No. 6, pp. 113–128 (2005).
- [21] 平尾 努, 鈴木 潤, 磯崎秀樹: 識別学習による組合せ最適化問題としての文短縮手法, *人工知能学会論文誌*, Vol. 22, No. 6, pp. 574–584 (2007).
- [22] 平尾 努, 鈴木 潤, 磯崎秀樹: 構文情報に依存しない文短縮手法, *情報処理学会論文誌: データベース*, Vol. 2, No. 1, pp. 1–9 (2009).
- [23] 長谷川隆明, 西川 仁, 今村賢治, 菊井玄一郎, 奥村 学: 携帯端末のための Web ページからの概要文生成, *人工知能学会論文誌*, Vol. 25, No. 1, pp. 133–143 (2010).