

# テキスト自動要約に関する最近の話題

奥村 学<sup>†</sup> 難波 英嗣<sup>†</sup>

本稿では、1999年の解説の後を受け、テキスト自動要約に関する、その後の研究動向を概観する。本稿では、その後の動向として、特に最近注目を集めている、以下の3つの話題を中心に紹介する。

- (1) 単一テキストを対象にした要約における、より自然な要約作成に向けての動き、
- (2) 複数テキストを対象にした要約研究のさらなる活発化、
- (3) 要約研究における、要約対象の幅の広がり

キーワード: テキスト自動要約, 首尾一貫性, 読み易さ, 冗長性, 自由作成要約, テキストのジャンル, 要約の書き換え, 要約の言語モデル, 複数テキスト要約, 話し言葉の要約

## New Topics on Automated Text Summarization

MANABU OKUMURA<sup>†</sup> and HIDETSUGU NANBA<sup>†</sup>

In this article, we try to survey the current trends in the field of automated text summarization, especially concentrating on the following three topics: researches on producing more natural summaries in single document summarization, farther activation of researches on multi-document summarization, and more variety of summarization inputs in the researches.

**Keywords:** *automated text summarization, coherence, readability, redundancy, human-written summaries, genre of texts, revision, language modeling, multi-document summarization, speech summarization*

### 1 はじめに

電子化されたテキストが世の中に満ち溢れる現状から、テキスト自動要約研究が急速に活発になり、数年が早くも経過している。研究の活発さは依然変わらず、昨年も NAACL に併設する形で要約に関するワークショップが6月に開催された。また、日本では、国立情報学研究所の主催する評価型ワークショップ NTCIR-2 のサブタスクの1つとしてテキスト自動要約 (TSC: Text Summarization Challenge) が企画され、日本語テキストの要約に関する初めての評価として、また、Tipster における SUMMAC に続く要約の評価として関心を集め、昨年3月にその第1回 (TSC1) の成果報告会が開催された (<http://research.nii.ac.jp/ntcir/index-ja.html>)。一方、アメリカでは、SUMMAC に続く評価プログラムとして、DUC (Document Understanding

<sup>†</sup> 東京工業大学精密工学研究所, Precision and Intelligence Laboratory, Tokyo Institute of Technology

Conference) が始まり、第1回の本格的な評価が昨年夏行なわれ、9月に開催された SIGIR に併設する形でワークショップが開催された (<http://www-nlpir.nist.gov/projects/duc/>).

このような背景の元、本稿では、1999年の解説(奥村, 難波 1999)の後を受け、テキスト自動要約に関する、その後の研究動向を概観する。1999年の解説では、これまでのテキスト自動要約手法として、重要文(段落)抽出を中心に解説するとともに、当時自動要約に関する研究で注目を集めつつあった、いくつかの話題として、「抽象化、言い換えによる要約」、「ユーザに適応した要約」、「複数テキストを対象にした要約」、「文中の重要箇所抽出による要約」、「要約の表示方法」について述べている。本稿では、その後の動向として、特に最近注目を集めている、以下の3つの話題を中心に紹介する。

- (1) 単一テキストを対象にした要約における、より自然な要約作成に向けての動き、
- (2) 複数テキストを対象にした要約研究のさらなる活発化、
- (3) 要約研究における、要約対象の幅の広がり

(1)の動きは、後述するように、1999年の解説における「抽象化、言い換えによる要約」、「文中の重要箇所抽出による要約」という話題の延長線上にあると言える。以下、2, 3, 4節でそれぞれの話題について述べる。

なお、TSC1 および DUC2001 にはそれぞれ多数の参加があり、興味深い研究も多い。しかし、TSC1 の多くの研究は重要文抽出に基づくものであり、本稿に含めるのは適当でないと考えた。また、DUC2001 に関しては、ワークショップが開催されたのが9月13, 14日であり、本稿に含めるのは時間的余裕がなく断念せざるを得なかった。これらについては、稿を改めて、概観することとしたい。

## 2 より自然な要約作成に向けて

ここ1, 2年テキスト自動要約研究者が関心を持っている話題に、単一テキストを対象にした要約において、人間にとってより自然な要約を目指すというものがある。

これまでの要約手法である重要文抽出には、問題点として、テキスト中の色々な箇所から抽出したものを単に集めているため、抽出した複数の文間のつながり(首尾一貫性)が悪いことが指摘されている。抽出した文中に指示詞が含まれていても、その先行詞が要約中に存在しない可能性があったり、また、不要な接続詞があったりするということだが、こういうことが起きると、読みにくいということはもちろんだが、最悪の場合、要約テキストの内容を読み間違えてしまう可能性もある。また、文を重要として要約に含める際、他の文とは独立に抽出を行っており、そのため、結果として要約中に抽出された文の内容に類似のものがいくつも含まれるということが生じる可能性がある。

このような、これまでの要約手法の問題点を受けて、「より読み易い要約」、「より冗長性の少ない要約」を目指す動きが近年活発になっており、また、人間の自由作成要約(human-written

summary) を元に要約手法を検討する動きも盛んになってきている。

人間が自由に要約を作成する際、原文に基づかず一から要約を「書く」場合もあるが、多くの場合、原文を元に、原文の断片を適切に「切り貼り」し、その後それに編集を加えることで、要約を作成しているという観察を元に、そういった人間の要約作成過程を計算機上にモデル化しようという研究も、後述するように(2.2節)始まっている(Jing and McKeown 2000)。人間の要約作成モデルに基づく要約手法なら、人間の要約に(ある程度)近い要約を作成できる可能性があり、注目すべき研究と言える。

もう一つ特筆すべき研究として、自然言語生成システムを利用した要約手法の提案も始まっている(McKeown, Klavans, Hatzivassiloglou, Barzilay, and Eskin 1999; Barzilay, McKeown, and Elhadad 1999)。詳細は3節で述べるが、複数テキスト中の重要箇所を、FUF/SURGEという生成システムにより、つなぎ合わせることで要約として生成している。

要約の過程は、大きくテキストの解釈(文の解析とテキストの解析結果の生成)と(テキスト解析結果中の重要部分の)要約としての生成に分けられるとされてきたが、これまでの研究では、要約を生成するという事は実際にはほとんど実現されていなかった。今後、より自然な要約作成を目指す過程で、自然言語生成技術の利用は不可欠となっていくであろう。

これまでも、要約の読みにくさ、首尾一貫性の悪さに対しては、対処法が提案されてきているが(たとえば、Mathisら(Mathis, Rush, and Young 1973)や(奥村, 難波 1999)の2.3節を参照)、いずれも ad hoc な手法という印象が強い。これに対して、抽出した重要文集合を書き換える(revise)ことで、文間のつながりの悪さを改善し、より読み易い要約作成を目指す研究が最近試みられている(難波, 奥村 1999)。まだ技術的に難しい問題がいろいろあるが、興味深い。

また、重要文抽出ではなく、文中の重要箇所抽出、不要箇所削除による要約手法はすでに(奥村, 難波 1999)で紹介されているが、この要約手法も、より自然な要約を作成するための第一歩と言える。2.4節で紹介する「要約の言語モデル」は、この要約手法を統計的に定式化した枠組とも考えられる。

以下、各小節で、「より冗長性の少ない要約作成」、「人間の自由作成要約を元にした要約手法」、「抽出した重要文集合の書き換えによる、より自然な要約作成」、「要約の言語モデル」の4つの話題について言及する。

## 2.1 冗長性の少ない要約に向けて

複数テキストを対象にした要約では、複数のテキストから抽出した内容を要約とする際、内容が重複することを避ける手法がとられることが一般的である。単一テキストを対象にした要約作成でも、要約中に類似した文が含まれていれば冗長であり、冗長性を削減することで、他の有用な情報を要約に加え、要約中の情報の密度を増すことができる。近年単一テキストの場合にも、要約中の冗長度を下げ、同じ長さの要約に、より多くの情報を含められるよう考慮し

た要約手法がいくつか提案されている。

Baldwin ら (Baldwin and Morton 1998) は、照応解析に基づき、query-sensitive で indicative(指示的) な要約を作成する手法を提案している。テキスト中の文を選択するのだが、検索要求中の句がすべて要約の中にカバーされるように選択する。テキスト中の句がその句と相互参照していれば、検索要求中の句はカバーされているとする。

文を選択する基準は、その文により新たにカバーされる (すでに選択された文ではカバーされていない) 検索要求中の句が多い文を選択する。この文選択をすべての句がカバーされるまで繰り返す。これにより、要約の冗長性を最小にしている。

Baldwin らの手法は、なるべく冗長な参照句を含まないように文を選択していることに相当する。また、先行詞を要約中に含まない代名詞は、可能なら先行詞に置き換える、不要と考えられる、前置詞句、同格の名詞句、関係節は除去するなどの後処理も施している。

MMR(Maximal Marginal Relevance)(Carbonell, Geng, and Goldstein 1997; Carbonell and Goldstein 1998) は、テキスト検索、単一テキスト要約、複数テキスト要約において利用可能な尺度であり、検索要求との適合度と、情報の新規性 (すでに選択されたものとの異なり度) をともに考慮する尺度である。MMR は、テキスト検索を例にすれば、以下の式で定義される。

$$MMR(Q, R, S) = \underset{D_i \in R \setminus S}{\operatorname{Argmax}} [\lambda \operatorname{Sim}_1(D_i, Q) - \\ (1 - \lambda) \max_{D_j \in S} \operatorname{Sim}_2(D_i, D_j)]$$

ここで、

$Q$ : 検索要求、

$R$ : システムによって検索された (ランク付けられた) テキスト群、

$S$ : すでに選択された  $R$  の部分集合

$R \setminus S$ :  $R$  と  $S$  の差集合

であり、 $\lambda$  は、検索要求との適合度 ( $\operatorname{Sim}_1(D_i, Q)$ ) と、すでに選択されたものとの異なり度<sup>1</sup>に関する重みづけ (どちらを重視するか) に関するパラメタである。なお、検索要求との適合度、すでに選択されたものとの類似度を計算する尺度  $\operatorname{Sim}_1, \operatorname{Sim}_2$  には、任意のものが利用できるが、単語を要素とするベクトル間の距離尺度 (たとえば、コサイン、内積等) を利用することが多い。

MMR を用いた要約では、query-relevant な要約を作成するが、単一テキスト要約では、検索要求に関連するパッセージ (文) の集合を ( $\operatorname{Sim}_1$  のみを利用して) まず抽出した後 (これが  $R$ )、それらを MMR で再順序付け、要約の長さまで文を選択し、原文での順序、MMR のスコアの順序等を元に出力する。したがって、1 文目は、検索要求と最も適合する文が選択され、2 文目以後は、それまでに選択された文 (2 文目の場合は、最初に選択された文) との異なり度も合わせて考慮して選択される。MMR を用いることで、要約は互いに (最大限) 異なる文により構成され

<sup>1</sup>  $\max_{D_j \in S} \operatorname{Sim}_2(D_i, D_j)$  が類似度を表しているのだから、それを引くことで、異なり度としている。

る。MMRを用いた複数テキスト要約は3節で紹介する。

加藤ら(加藤，浦谷 2000)は，放送ニュースを対象にした重要文抽出法として，まず1文目(リード文)を抽出した後，それ以後の文のうち，リード文と内容が重複しない文を重要として抽出する手法を提案している。内容の重複は，文間の単語の対応の度合を元に計算している。この手法は，重要文抽出に，テキスト中での位置情報とMMRの考え方を併用していると言うことができる。

石ごころ(石ごこ，片岡，増山，中川 1999)は，同一の事象を表す表現が複数回テキスト中に出現した場合，2回目以後の出現を重複部分として削除する手法を提案している。

## 2.2 人間の自由作成要約を目指して

人間は，単に重要文を抽出するだけでなく，それらを編集することで要約を作成していると考えられる。Jingら(Jing and McKeown 1999, 2000)は，人間の自由作成要約と原文の対応を分析し，抽出された文を編集する6つの操作を同定している。それらは，不要な句の削除(文短縮)，(短縮した)文を他の文と結合する(文の結合)，構文的変形，句を言い替える(語彙的言い替え)，句をより抽象的/具体的な記述に置き換える，抽出した文を並べ替える，の6つである。一方，人間が原文に基づかず，一から書いている文も自由作成要約には含まれており，その割合は，300要約を調べたところ，19%であったと報告している。

Jingら(Jing and McKeown 2000)は，人間の自由作成要約の分析から得られた6つの編集操作を用いた「切り貼り」に基づく要約手法を提案している。システムは，抽出された重要文を編集し，不要な句を削除し，結果として残った句をまとめ上げることで一貫性のある文を作成する。Jingらの切り貼りに基づく要約システムは，まず重要文を抽出した後，抽出した文を，6つの操作で(文短縮，文の結合のみが実装されている)編集し，その結果を要約として出力する。文の結合に関しては，対応コーパスを分析し，人手で規則を作成して実現している。文の結合は，2つの構文解析木に対する，結合，部分木の置換，ノードの追加というTAG上の操作として実装されている。

一方，文短縮は，抽出された重要文から，不要な句を自動的に削除するが，人間の自由作成要約と原文の対応コーパスから得られた統計情報，構文的知識，文脈情報を利用して，削除する句を決定している(Jing 2000)。

原文は，構文解析され，構文解析木中の必須要素と考えられる部分は印が付けられ，後の処理で削除され，文法的でない文が作成されることを防止する。次に，文中の句で話題ともっとも関連するものを決定する。また，対応コーパスを構文解析した結果を用いて，どの句がどういふ条件でどの程度削除され易いか(たとえば，主動詞が‘give’のとき，‘when’節が削除される確率)を計算する。また，句が短縮される(部分が削除される)確率，句が変化しない確率も合わせて計算される。そして，必須でなく，話題とあまり関係がなく，人間が削除している確

率がある程度ある句を削除の対象とする。

人間の削除箇所との一致度に基づく評価では、平均で 81.3%の精度を得ており、すべての前置詞句、節、to 不定詞、動名詞を削除する場合を baseline と考えるなら、baseline の精度は 43.2%だった。また、システムは平均で文の長さを 32.7%短くしていたが、人間の場合は 41.8%だった。システムの出力における誤りの原因は、50 文を分析した結果では、8%が構文解析誤りによるものだった。

この Jing らの研究と同様、(重要文抽出ではなく、) 人間が自由に作成した要約のコーパスに基づいた要約研究が近年数多く見られる。これらの研究では、人間の自由作成要約と原文を対応付けた (aligned) コーパスが必要であるため、要約と原文の間の対応づけ (alignment) を行なう手法に関する提案もいくつか見られる。

Jing ら (Jing and McKeown 1999) の対応づけプログラムは、人間の自由作成要約中の句を原文中の句に自動的に対応付ける。要約中で隣接する 2 単語は、原文中でも隣接して現れ易い、遠く離れた文中に現れないというようなヒューリスティクスを元にした HMM に基づいており、要約中の各単語が原文中のどこに位置するかを Viterbi アルゴリズムにより決定する。50 要約中の 305 文に対する対応関係を人手で調査したところ、93.8%の文で正しい対応関係を得ていると報告している。

Marcu (Marcu 1999) は、原文と自由作成要約をとともに、出現する単語のベクトルで表現し、その間の類似度をコサイン距離で計算する。そして、自由作成要約と類似度がもっとも大きくなるように、原文から節を削除していくことで、対応する抜粋を決定している。

Banko ら (Banko, Mittal, Kantrowitz, and Goldstein 1999) は、文を単位とし、文を文中の単語の出現頻度のベクトルで表し、ベクトル間の距離で文間の類似度を計ることで、自由作成要約中の文と原文中の文をもっとも類似度が大きくなるように対応付けている。Banko らと Marcu の手法はともに、abstract から抜粋 (extract) を生成することを目的としているため、対応させる単位が文、節と大きい。

望主ら (望主, 荻野, 太田, 井佐原 2000) も、自由作成要約を原文と対応付けるツールを作成し、対応結果から、自由作成要約、重要文抽出による要約の相違点の分析を行なっている。また、(奥村, 難波 1999) で紹介されている加藤らは、要約知識の自動獲得を目的に、単語の部分一致を考慮した DP マッチングによる対応づけ手法を示している。

このようにして、自由作成要約と原文を対応付ける (あるいは、対応する抜粋を生成する) と、自由作成要約と抜粋の間の比較・分析が可能になる。

Marcu (Marcu 1999) は、人間の要約に含まれる内容をすべて含むように、テキストの抜粋を作成する場合、どの程度の長さの抜粋が必要であるかを調査している。新聞記事を対象にした場合、対応する要約と比べ、抜粋は 2.76 倍の長さが必要であるという結果を示している。この結果は、抜粋中の冗長性を除去したり、さらに文をより短くするなど、抜粋をさらに加工す

る必要があることを示しているとも言える。

また、Jingらは、自由作成要約は、対応する抜粋と比較すると、52%の長さであるという報告をしている。Goldsteinら (Goldstein, Kantrowitz, Mittal, and Carbonell 1999) の報告では、平均して抜粋の長さは、自由作成要約に比べ、20%長くなるという。

## 2.3 要約における言い替え、書き換えの役割

2.2節で述べたように、人間の要約過程は、単に重要文を抽出するだけでなく、それらを編集する操作が含まれていると考えられる。この編集の操作には、書き換え (revision) や言い替え (paraphrase) が含まれている。本節では、書き換えや言い替えが用いられた要約研究を概観する<sup>2</sup>。

抽出した重要文集合である抜粋を書き換える目的には、少なくとも次の2つがあると考えられる。

- (1) 文の長さを短くする
- (2) 抜粋を読み易くする

片岡ら (片岡, 増山, 山本 1999) は、連体修飾節を含む名詞句を「AのB」の形に言い替えることで要約を行なう手法を示しているが、これは前者に該当すると言える。また、(奥村, 難波 1999) で紹介されている、概念辞書等を用いて語句を抽象化する言い替えを行ない要約する手法である「抽象化, 言い換えによる要約手法」(3節)や、加藤, 若尾らのような手法(6節)は、言い替えを行なうことで、文字列を削減する要約手法とすることができる。

また、Maniら (Mani, Gates, and Bloedorn 1999) は、抜粋を書き換えることで、質の向上を目指している。3つの操作, elimination, aggregation, smoothingを示している。それらを抜粋に繰り返し適用することで、抜粋の読み易さを低下させずに informativeness を向上できたと主張している。このことから、Maniらの主眼は、書き換えにより、要約内の情報の量を向上させること(抜粋中の不要な個所を削除することで、他の個所の情報を要約に加える)であると言える。

eliminationがJingらの文短縮, aggregationとsmoothingが文の結合にそれぞれ対応している。eliminationでは、文頭の前置詞句, 副詞句を削除する。smoothingには、読み易さ(首尾一貫性)を改善するための操作が一部含まれる。

一方、後者の研究としては、難波ら(難波, 奥村 1999)の研究がある。難波らは、人間に抜粋を書き換えてもらう心理実験を行ない、抜粋の読みにくさの要因を分析した後、要因ごとに読みにくさを解消するための書き換えを定式化している。接続詞を追加したり、削除したり、また、冗長な単語の繰り返しが代名詞化したり、省略したり、逆に、省略されている単語を補完したり、などである。そして、そのうちいくつかを実装している。

2 川原(川原 1989)は、人間の要約作成過程において、どのように書き換えが役割を果たしているかを調査している。

大塚ら(大塚, 内海, 奥村 2001)は,「この」等の指示形容詞を含む名詞句に対して照応処理を行なうことで,対応する先行名詞句を特定し,指示形容詞を含む名詞句に対応する先行名詞句に置き換えることで,抽出した重要文集合のつながりの悪さを改善する手法を示している。

## 2.4 要約の言語モデル

原文と自由作成要約の組がコーパスとして大量に存在するならば,人間の要約過程を模倣するようにモデルを訓練することが可能である。Knight と Marcu(Knight and Marcu 2000)は,このような考え方にに基づき,文要約(文短縮)において,文法的で,しかも,内容としては原文の情報の重要な部分を維持するような手法を2つ示している。2つの手法は,確率的 noisy-channel モデルと決定木をそれぞれ用いている。入力として,単語列(1文)を与えると,単語列中の単語の部分集合を削除し,残った単語が要約を構成する<sup>3</sup>。

確率的 noisy-channel モデルは,統計的機械翻訳の場合と同様,次の2つのモデルで構成される。

- Source Model:

要約を構成する文  $s$  の確率  $P(s)$ 。文  $s$  が生成される確率を示す。この確率は,文法的でない文の場合低くなり,要約が文法的であるかどうかの指標となる。単純には bigram でモデル化される。

- Channel Model(Translation model):

単語列の組  $\langle s, t \rangle$  の確率  $P(t|s)$ 。要約  $s$  から,より長い単語列  $t$ (原文)が得られる確率。原文中の各単語が要約に出てくる確からしさを示しており,各単語の確からしさの積をその単語列が要約となる確からしさとする。重要な内容を保持しているかどうかの指標となる。

Knight と Marcu は,上の2つの確率を単語列に対してではなく,それを構文解析した結果得られる木に対して計算している。 $P_{tree}(s)$  は,木  $s$  を得る際に利用される文法規則に対して計算される標準的な確率文脈自由文法のスコアと,木の葉に現れる単語に対して計算される標準的な単語の bigram のスコアの組合せである。

確率的な channel モデルでは,拡張テンプレートを確率的に選択する。たとえば, NP と VP を子ノードとして持つノード  $S$  に対して,確率  $P(S \rightarrow NPVP PP|S \rightarrow NPVP)$  を元に,子ノード PP を追加する。

そして,単語列  $t$  からそれに対応する要約  $s$  を選択する際,  $P(s|t)$  を最大にするものを選択する。これは,  $P(s) \times P(t|s)$  を最大にする  $s$  を選択することと同じである。原文中の単語列の部分集合で,上の2つの確率の積を最大にするものを Viterbi ビームサーチを用いて選択する。

<sup>3</sup> 原文と抜粋の組のコーパスから重要文抽出のためのモデルを学習する手法については(奥村, 難波 1999)の2.2節すでに紹介されている。



Ziff-Davis コーパス中の 1067 組の文を対象にし訓練を行なっている。拡張テンプレートは、原文と要約文とともに構文解析し、その木の対応関係から抽出している。

一方、決定木に基づく手法としては、原文に対応する木  $t$  を与えると、それを要約文に対応する、より小さな木  $s$  に書き換えるモデルを示している。拡張した決定的 shift-reduce 構文解析の枠組に基づき、空のスタックと、入力の木  $t$  を入れた入力リストを用いて処理を開始し、より小さな木へ書き換えるべく、shift (入力リストの先頭をスタックへ移動), reduce (スタック上の  $k$  個の木を組み合わせて新たな木を構成し、スタックにプッシュ), drop (入力リスト中の構成素を削除) の操作を繰り返し実行する。

決定木に基づく手法は、noisy-channel モデルに基づく手法よりも、より柔軟であり、原文の構造と要約文の構造が著しく異なる場合にも対処可能である。どの操作を選択するかは、訓練データ (原文-要約文の組の集合から構成される操作の系列の集合) から、決定木学習を行なうことで学習される。

このように、文要約のモデルを、訓練コーパスから自動的に訓練することで得る手法は、Witbrock と Mittal (Witbrock and Mittal 1999) が、原文と abstract の組で直接訓練した確率モデルを適用したのが最初の研究とされる。これ以外は、前節で紹介した Jing らの研究や、(奥村, 難波 1999) で紹介されている、文中の重要個所抽出、不要個所削除による要約手法を含め、いずれも、人手で作成した、あるいは半自動で得た規則を元に、冗長な情報を削除したり、長い文をより短い文に縮めたり、複数の文をまとめたりしている。

堀之内ら (堀之内, 山本 2000) は、「日本語らしく、かつ意味的に重要個所を含む」ように、文を短縮する統計的手法を示している。日本語らしさの評価のために n-gram モデル、意味的に重要個所を含むかどうかの評価のために idf をそれぞれ利用している。この 2 つを重み付けした重要度を文中の断片に与え、重要度の小さい断片を繰り返し削除することで文を短縮していく。

小堀ら (小堀, 田村 2000) は、あらかじめ原文から抽出された重要文節データを元に学習した決定木を用いて重要文節を抽出する手法を示している。

Berger と Mittal (Berger and Mittal 2000b) は、query-relevant な要約を作成する統計的言語モデルを示している。FAQ のコーパスを訓練データとして、文書  $d$  とクエリ  $q$  の組に対して、

$$p(s|d, q) = p(q|s, d) * p(s|d) \approx p(q|s) * p(s|d)$$

を最大にする  $s$  を要約として求める。そして、そのための確率  $p(q|s), p(s|d)$  をそれぞれ訓練データから学習する。確率  $p(q|s), p(s|d)$  はそれぞれ、(クエリに対する要約の) 適切性 (relevance)、(テキストに対する要約の) 忠実性 (fidelity) と呼ばれている。

### 3 複数テキストを対象にした要約手法

これまでの複数テキスト要約研究では、あらかじめ人間が用意した比較的小規模なテキスト集合をシステムの入力として要約を作成するのが中心的であったと言える。しかし、近年、情報検索システムの検索結果を直接要約システムの入力に用いるなど、より大規模なテキスト集合を要約対象とする実用性の高いシステムがいくつか提案されてきている。

要約システムの入力として想定されるテキスト集合は、(1) すべてが同一トピックのものと、(2) 情報検索システムの検索結果のように、複数のトピックが混在しているものの大きく2種類存在すると考えられる。どちらのテキスト集合を対象とするかで、要約作成手法、要約システムの位置付けも次のように異なってくる。

- (1) 要約システムに与えるテキスト集合中のテキストはどれも同じトピックについて書かれたものであり、そのため、似たような内容のテキストが複数含まれる可能性がある。この場合、すべてのテキストの内容を要約に含めると、冗長な要約が作成されてしまう。そこで、テキスト(あるいはテキスト中のパッセージ)間の類似度を考慮し、内容がなるべく重複しないように要約を作成する。
- (2) 情報検索の結果得られたテキスト集合を要約システムの入力に用いるような場合、そのテキスト集合には、ユーザの目的と合致しないテキストが数多く含まれている可能性がある。このような場合、目的のテキスト集合へユーザをナビゲートする支援システムは有用であり、そのようなシステムでは、テキスト集合を自動的に分類し、グループごとに、グループのテキスト集合の要約を作成しラベルとして付与する。ユーザは、自分の必要なテキストがグループに含まれているかどうかを付与されたラベルを見て判断する。

(奥村, 難波 1999)では複数テキスト要約のポイントとして、図1に示す3点を挙げて、この3点に沿って研究を概観していた。本節でも同様にこの3点に沿って、この分野の最近の研究動向を、上の分類に即して、紹介する。

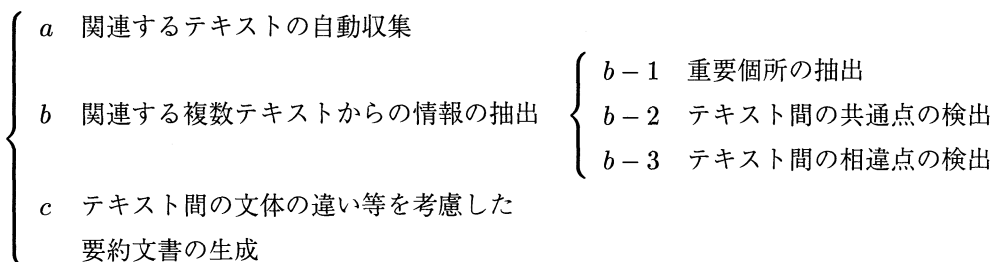


図1 複数テキスト要約のポイント

## 分類 1: 同一トピックのテキスト集合からの要約作成

分類 1 の要約手法には Goldstein ら (Goldstein, Mittal, Carbonell, and Kantrowitz 2000), Radev ら (Radev, Jing, and Budzikowska 2000), Stein ら (Stein, Strzalkowski, and Wise 1999), Barzilay ら (Barzilay et al. 1999; Barzilay, Elhadad, and McKeown 2001), McKeown ら (McKeown et al. 1999) のものがある。

Goldstein ら (Goldstein et al. 2000) は新聞記事を対象とし、記事集合中からある検索クエリに関するパッセージを抽出、収集し (a), それらを並べて要約を作成する MMR-MD (Maximal Marginal Relevance Multi-Document) という手法を提案している。検索されたパッセージを単純にクエリとの適合度の高い順に並べただけでは、パッセージ間で重複する個所が存在する可能性があり、要約として望ましくない。そこで、MMR-MD では、クエリに対するパッセージの適合度を考慮しつつ、すでに上位にランクされているパッセージと類似度の低いもの (重複個所が少ないと思われるパッセージ) (b-3) を選択して順に出力することで、冗長性の少ない複数テキスト要約の作成を行っている。また、パッセージの出力順序を決める際、記事が書かれた日時なども考慮している。

Radev ら (Radev et al. 2000) は新聞記事集合をあらかじめクラスタリングし、各クラスタごとに要約を作成する手法を提案している (a)。クラスタ中の記事中の各文の重要度をまず計算し、次に要約率に応じて記事集合から重要度の高い文を抜き出し、抜き出された文を記事の書かれた日付順に並べて、要約として出力する。文の重要度は、クラスタの特徴を表す語を文が含む割合 (b-2), 文の位置 (lead) (b-1) により決定する。また、Goldstein らと同様、自分より重要度の高い文と内容が重複するような文は重要度を下げること、冗長性の少ない要約の作成を目指している (b-3)。

Stein ら (Stein et al. 1999) は、あらかじめテキストごとの要約を作成し (b-1), 作成された要約をクラスタリングし、似たような内容の要約をグルーピングしている。そして、各クラスタ中で最も代表的な要約をクラスタの要約として抽出する (b-2)。また、クラスタの要約同士の類似度を計算し、隣接する 2 つの要約の類似度が高くなるよう並べ換えて出力している。

Barzilay ら (Barzilay et al. 1999), McKeown ら (McKeown et al. 1999) は、複数の新聞記事間で言い回しは異なるが同じ内容の文を、7 種類の言い換え規則を用いて同定している (b-2)。同定された文は、構文解析器を用いて述語項構造に変換され、文間で共通な句が抽出される。その後、文生成器を用いて抽出された共通語句を統合し、要約文として出力する (c)。さらにこれらの要約文は記事の日付順およびテキスト中の出現順にソートされ、それらが最終的な要約文書となる。

要約文書の構成要素となるトピック (文) を並べる順序を決定するこれまでの方法は、文間のつながりを考慮する方法 (Goldstein et al. 2000; Stein et al. 1999) と、記事が書かれた時間順に並べる方法 (Radev et al. 2000; McKeown et al. 1999) の 2 つに分けられる。一般にはさまざま

まなトピックの並べ方が存在するが、人間が複数のテキストから要約を作成する場合、何らかの原理に基づいて並べる順序を決定していると考えられる。Barzilay ら (Barzilay et al. 2001) は、複数の記事から抽出されたいくつかの重要文のセットを 10 人の被験者に与え、それらを並べ換えることで要約を作成してもらっている。そして、その結果を比較することで、次のような知見を得ている。

すべての文の順序が被験者間で完全に一致することはあまりない。しかし、順序が入れ替わっても、常に隣り合って出現する文のペアがいくつかある。これらのペアは関連したトピックの文で構成されている。したがって、複数テキスト要約において文間の結束性を考慮することは重要である。

このような知見に基づき、Barzilay ら (Barzilay et al. 2001) は、要約文の順序を決定する方法を考案している。基本的にはトピックを時間順に並べるが、関連したトピックの文は必ず隣接して出力する。この方法により、作成される要約文書はある程度結束性が保たれる。Barzilay らは、この手法を先に述べた「記事が書かれた時間順に並べる方法」と比較し、前者の手法の方が優れていることを示している。

## 分類 2: 複数のトピックを含んだテキスト集合からの要約作成

分類 2 の要約手法には Eguchi ら (Eguchi, Ito, Kumamoto, and Kanata 1999), Fukuhara ら (Fukuhara, Takeda, and Nishida 1999), Ando ら (Ando, Boguraev, Byrd, and Neff 2000), 上田ら (上田, 小山 2000), Kan ら (Kan, McKeown, and Klavans 2001) のものがある。

Eguchi ら (Eguchi et al. 1999) は、WWW 上のテキストを対象にした関連性フィードバックに基づく検索システムを構築している。このシステムでは、検索結果 (a) をテキスト間の類似度に基づいてクラスタリングし、各クラスタごとにクラスタに多く含まれる語と、そのクラスタを代表するテキストのタイトルを、そのクラスタの要約として出力する (b-2)。出力されたクラスタをユーザに選択してもらい、そのクラスタに含まれるテキストを用いて関連性フィードバックを行っている。

Fukuhara ら (Fukuhara et al. 1999) も、Eguchi らと同様に検索結果をクラスタリングし (a)、クラスタごとに要約出力を行っている。Fukuhara らは、テキスト中の単語の出現頻度分布を考慮し、クラスタごとの話題を表す語とそれらを含んだ文を抽出する。さらに、抽出された文を、焦点-主題連鎖を考慮して並べ替え、クラスタごとの要約として出力している (b-2)。

Ando ら (Ando et al. 2000) は、ベクトル空間モデルを用いて新聞記事集合中の記事間の類似度を計算し、それらを semantic space と呼ばれる 2 次元空間上に配置し表示するシステムを構築している。semantic space 上では各記事はドットで表現され、またトピックの似た記事は semantic space 上で隣接して配置される。マウスで semantic space 上のドットを指せば、そのドット (記事) と関連のあるドット (記事) が強調され (a)、さらに関連記事中の頻出単語 (topic

term) や頻出単語を多く含む文 (topic sentence)(b-2) が表示される。

上田ら (上田, 小山 2000) は, クラスタリングによりある程度同じ話題でまとめられたテキスト集合を対象に, 各クラスタの特徴を表す文を自動的に作成する手法を提案している (a). 上田らも Barzilay ら, McKeown らと同様に, テキスト中の各文を構文解析し, テキスト間で構文木同士を比較することで, テキスト間の共通個所を同定するという手法を提案している (b-2). 構文木の比較には 2 種類の方法を提案している. 1 つは, 例えば「フーバー社が携帯電話を発売」という文を, 意味的に等価な「携帯電話がフーバー社から発売」などに構文レベルで変換し, 同一内容の異なる 2 文を同定し, クラスタのラベルとして出力するという方法である. もう 1 つは, シソーラスを用いて「ホウレンソウからダイオキシシンが検出された」と「白菜からダイオキシシンが検出された」の 2 文から「野菜からダイオキシシンが検出された」といったように, より抽象度の高いレベルで融合し, ラベルとして出力するという方法である.

これまで分類 2 で述べてきた要約作成手法は, 形態素, あるいは構文レベルでテキスト間の比較を行っているが, 個々のテキストからいくつかの属性値を抽出した後, テキストを属性レベルで比較し, 要約を作成する試みがある.

Kan ら (Kan et al. 2001) は, ある検索クエリで検索された医療関係のテキスト集合を比較し, ユーザがどのテキストを読むべきか判断するのに有用な indicative(指示的) な要約を生成する手法を提案している. このシステムを利用することで, 例えば喉頭炎 (angina) を検索クエリとした場合, 「検索結果は 23 件あります. 結果には喉頭炎のガイド ('the AMA Guide to Angina') が含まれています.」「喉頭炎に関する定義やリスクに関して述べたテキストがあります.」「喉頭炎の関連情報を含んだテキストがあります.」といった要約が出力される. このような要約を生成するために, Kan らは, まず個々のテキストを, セクションの情報に基づいて, そのテキストのトピックの構造を示す木 (トピックツリー) で表現する. 次に検索クエリがトピックツリーのどこに位置するのか (クエリとテキストのメイントピックとの関連), テキストの平均的な属性値と比較して, テキストに含まれるいくつかのトピックがどの程度重要であるのか (他のテキストと異なっているのか) (b-1, b-3) といった情報をテキストの属性値として抽出する. これらの値を用いて要約生成を行い, 検索結果として出力する.

## その他の要約作成手法

これまでの複数テキスト要約の研究は, 複数のテキストから得られた情報をいかに統合して要約を作成するかに主眼が置かれてきたと言える. 一方, ある種のテキスト集合には, 集合全体の内容をまとめたテキスト (パッセージ) が存在することがある. 例えば, ある分野の研究動向をまとめたサーベイ論文や, ある事件に関する解説記事などがこれに相当する. このようなテキストを見つけ出すことができれば, それ自体を複数テキストの要約とみなすことができる.

橋本ら (橋本, 奥村, 島津 2001) は, ある事件に関して過去の主要な出来事が新聞記者の観

点で要約されている個所をサマリパッセージと呼び、このような個所を記事集合から自動的に抽出する手法を提案している。サマリパッセージは、解説記事や社説など何らかの意見が述べられている記事(意見記事)中に含まれている。また事件の経緯を時系列に箇条書でまとめた個所もサマリパッセージと考えることができる。橋本らは、表層的な情報を用いてこれらのサマリパッセージの抽出を試みている。例えば、記事のタイトルに「社説」や「解説」を含む記事、意見文を多く含む記事は、意見、解説記事として抽出される。また、新聞固有の箇条書の形式を認定することでまとめ記事が抽出できる。こうして抽出されたパッセージは人間が作成したものであるため、これまでの複数テキスト要約で問題とされてきた要約の一貫性が保証されている。

## 4 要約対象の幅の広がり

これまでの自動要約研究の多くは、その要約対象のテキストのジャンルとして、新聞記事、論文を扱ってきた。これに対し、近年これ以外のジャンルのテキストを要約対象とする研究が見られるようになってきた。たとえば、web page を対象とした研究としては OCELOT(Berger and Mittal 2000a) 等があり、また、mail を対象とした研究としては(遠山, 西田 2000) 等がある。

さらに、テキストではなく、音声(あるいは、その書き起こしである話し言葉のデータ)を対象とする要約研究がいくつか見られるようになってきた。これには、講演音声のような monologue と、2人以上による対話(dialogue)の両方が含まれる。話し言葉を対象とした要約では、(1) テキストとしての情報以外に他の音響的信息が利用できる、(2) 音声認識結果を入力とすることから、入力にノイズが含まれる、(3) 後述するように、話し言葉の特性としての冗長性が入力には含まれる等、テキストを対象とした場合とは異なり、新たに考慮しなければいけない点が存在する。そこで、本節では、以下、これらの話し言葉を対象とした要約研究を概観する。なお、これまでも、(奥村, 難波 1999) で紹介されている字幕作成における要約のように、入力としてニュース原稿の読み上げ音声を対象とした研究は存在する。

堀と古井(堀智, 古井 2001) は、講演音声を自動要約する手法として、各発話文から重要な単語を抜き出し、それらを接合することで要約文を作成する手法を提案している。要約は、要約のもっともらしさを示す要約スコアを最大にする文中の部分単語列を DP マッチングにより決定し得ている。要約スコアは、単語の重要度(頻度に基づく)、単語連鎖の言語スコア(単語の trigram)、音声認識時の各単語の音響的、言語的信頼度、および原文中の単語の係り受け構造に基づく単語間遷移確率の重みつき和として定義される。

講演は、自然な発話(spontaneous speech)に比べれば整っているが、フィラーや言い直しなど、多くの冗長表現を含み、話し言葉に近い特性をもつ。この特徴を利用し、幅田と奥村(幅田, 奥村 2001) は、冗長表現を不要個所として削除することで、情報を欠落させずに要約を行

う手法を示している。人手によって講演音声の要約を行っている要約筆記データの分析をまず行い、その分析結果を元に、文短縮型の要約システムを開発している。分析の結果、フィルター、言い直し・繰り返し表現、挿入句表現、丁寧表現、「～という～」表現が削除または言い替えの対象として得られている。削除率および、要約筆記データを正解データとした場合の精度を尺度として要約システムを評価したところ、削除率 18.0%、精度 79.8%が得られている。この研究は、聴覚障害者のための情報保証手段の一つとして人手で現在行なわれている要約筆記の自動化を目指すものと言うことができる。

笠原と山下(笠原, 山下 2001)は、講演音声を対象とした要約の自動作成のため、重要文と韻律的特徴の関係についての分析を行なっている。

Zechner ら (Zechner and Lavie 2001) は、対話を書き起こしたものを入力とし、MMR により文をランク付けし、要約の長さまで、テキストの順序で文を出力する手法を示している。しかし、この手法では、質問に対応する応答が要約に含まれないため、一貫性に欠ける要約ができる可能性がある。そのため、複数の話者の発話にまたがる局所的な一貫性(この研究では、質問・応答の組のみ)を検出し、それを要約の際考慮に入れる(その一部がMMRで選択された場合組全体を要約に含める)ことで、要約の読み易さが向上することを人間の主観評価により示している。

Reithinger ら (Reithinger, Kipp, Engel, and Alexandersson 2000) は、音声翻訳システム VERBMOBIL を用いた、日程調整、ホテル予約のような領域における「交渉」対話を対象にした要約手法を示している。話し手の意図を発話行為クラスとして同定し、その情報を用いて、意図が suggest ならその内容を候補とし、reject なら棄却、give\_reason なら無視するというように、情報の選択の際に利用する。また、キーワードスポッティングにより、発話の内容を属性-値の組として同様に抽出する。そして、交渉対話では、話し手全員が合意したことに関心があるという前提を利用し、suggest された内容で、accept されたものを同定し、それを生成器で生成することで要約を作成している。

## 5 おわりに

1999年の解説(奥村, 難波 1999)の後を受け、テキスト自動要約の研究分野において、ここ数年関心が高まっている話題を3つ紹介した。

テキスト自動要約は、必要性が高まっていることもあり、今後も活発に研究が進められていくことと思われる。今後は、複数テキスト要約だけでなく、さらに対象範囲を広げ、複数の言語で書かれたテキスト(translingual summarization)、複数のメディアの情報を対象にした(テキストだけでなく、画像や音声も対象にする)要約(multi-media summarization)なども注目を集めそうである。今後も、テキスト自動要約の研究分野の動向には目が離せない。

また、テキスト自動要約技術の応用として、いくつかの新しい方向性が明確になってきた

ことも、ここ数年の話題と言えるかもしれない。これまでも、サーチエンジンにおける検索結果の表示や、ユーザのナビゲーションにおいて要約を利用する研究や、字幕作成、文字放送用に要約手法を利用することは試みられていた。これに加えて、ここ数年で、携帯端末における情報提示のための要約の利用(たとえば、(Buyukkokten, Garcia-Molina, and Paepcke 2001; Corston-Oliver 2001))や、(高齢者、視聴覚障害者といった)情報弱者のための情報保証への要約の利用(たとえば、自動要約筆記(幅田, 奥村 2001)やユーザの視覚特性に合わせたトランスコーディング(前田, 高木, 福田, 浅川 2001))といった、新しい有望な応用分野が要約には付け加わったと言える。研究分野の動向とともに、今後、要約の応用分野の動向にも目が離せないと言える。

最後に、新しい参考文献をいくつか紹介しておく。1999年に出版された(Mani and Maybury 1999)は、この分野の論文を、古典から最新のものまで集めた論文集であり、テキスト自動要約の最初の研究とも言われる(Luhn 1958)も入っている。この分野で研究を始める人には必読と言える。

TipsterのText Program Phase IIIの論文集(Tipster 1999)も出版されている。SUMMAC参加システムの概要がいくつか収録されており、また、SUMMACのdryrunの報告も含まれている。

また、昨年自動要約に関する教科書も出版されている(Mani 2001)。自動要約に関する話題をわかり易く記述しており、この本もこの分野で研究を始める人には必読と言える。なお、この本の翻訳の出版計画も進んでいる。

## 参考文献

- Ando, R., Boguraev, B., Byrd, R., and Neff, M. (2000). "Multi-Document Summarization by Visualizing Topical Content." In *Proc. of the ANLP/NAACL2000 Workshop on Automatic Summarization*, pp. 79–88.
- Baldwin, B. and Morton, T. (1998). "Dynamic Coreference-Based Summarization." In *Proc. of the 3rd Conference on Empirical Methods in Natural Language Processing*, pp. 1–6.
- Banko, M., Mittal, V., Kantrowitz, M., and Goldstein, J. (1999). "Generating Extraction-Based Summaries from Hand-Written Summaries by Aligning Text Spans." In *Proc. of the PACLING'99*, pp. 276–281.
- Barzilay, R., Elhadad, N., and McKeown, K. (2001). "Sentence Ordering in Multidocument Summarization." In *Proc. of HLT2001 1st International Conference on Human Language Technology Research*, pp. 149–155.
- Barzilay, R., McKeown, K., and Elhadad, M. (1999). "Information Fusion in the Context of Multi-Document Summarization." In *Proc. of the 37th Annual Meeting of the Associa-*



- tion for Computational Linguistics*, pp. 550–557.
- Berger, A. and Mittal, V. (2000a). “OCELOT: A system for summarizing web pages.” In *Proc. of the 23rd Annual International ACM-SIGIR Conference on Research and Development in Information Retrieval*, pp. 144–151.
- Berger, A. and Mittal, V. (2000b). “Query-Relevant Summarization using FAQs.” In *Proc. of the 38th Annual Meeting of the Association for Computational Linguistics*, pp. 294–301.
- Buyukkokten, O., Garcia-Molina, H., and Paepcke, A. (2001). “Text Summarization of Web pages on Handheld Devices.” In *Proc. of the NAACL2001 Workshop on Automatic Summarization*, pp. 109–110.
- Carbonell, J., Geng, Y., and Goldstein, J. (1997). “Automated Query-Relevant Summarization and Diversity-Based Reranking.” In *Proc. of the IJCAI-97 Workshop on AI in Digital Libraries*, pp. 9–14.
- Carbonell, J. and Goldstein, J. (1998). “The use of MMR, Diversity-Based Reranking for Reordering Documents and Producing Summaries.” In *Proc. of the 21st Annual International ACM-SIGIR Conference on Research and Development in Information Retrieval*, pp. 335–336.
- Corston-Oliver, S. (2001). “Text compaction for display on very small screens.” In *Proc. of the NAACL2001 Workshop on Automatic Summarization*, pp. 89–98.
- Eguchi, K., Ito, H., Kumamoto, A., and Kanata, Y. (1999). “Adaptive Query Expansion Based on Clustering Search Results.” *情報処理学会論文誌*, **40** (5), pp. 2439–2449.
- Fukuhara, T., Takeda, H., and Nishida, T. (1999). “Multiple-text Summarization for Collective Knowledge Formation.” In *Proc. of the Workshop on Social Aspects of Knowledge and Memory, IEEE Systems, Man and Cybernetics Conference*.
- Goldstein, J., Kantrowitz, M., Mittal, V., and Carbonell, J. (1999). “Summarizing Text Documents: Sentence Selection and Evaluation Metrics.” In *Proc. of the 22nd Annual International ACM-SIGIR Conference on Research and Development in Information Retrieval*, pp. 121–128.
- Goldstein, J., Mittal, V., Carbonell, J., and Kantrowitz, M. (2000). “Multi-Document Summarization by Sentence Extraction.” In *Proc. of the ANLP/NAACL2000 Workshop on Automatic Summarization*, pp. 40–48.
- 幅田隆, 奥村学 (2001). “不要個所削除による講演音声の要約.” *言語処理学会第7回年次大会発表論文集*, pp. 289–292.
- 橋本力, 奥村学, 島津明 (2001). “複数記事要約のためのサマリパッセージの抽出.” *言語処理学会第7回年次大会発表論文集*, pp. 285–288.

- 堀智織, 古井貞熙 (2001). “講演音声の自動要約の試み.” 話し言葉の科学と工学ワークショップ 講演予稿集, pp. 165–171.
- 堀之内寛, 山本幹雄 (2000). “n-gram モデルと IDF を利用した統計的日本語文短縮.” 言語処理学会第 6 回年次大会発表論文集, pp. 364–367.
- 石ざこ友子, 片岡明, 増山繁, 中川聖一 (1999). “テレビニュース番組の字幕作成のための重複部削除による要約.” 情報処理学会自然言語処理研究会報告, pp. 45–52. 133-7.
- Jing, H. (2000). “Sentence Reduction for Automatic Text Summarization.” In *Proc. of the 6th Conference on Applied Natural Language Processing*, pp. 310–315.
- Jing, H. and McKeown, K. (1999). “The Decomposition of Human-Written Summary Sentences.” In *Proc. of the 22nd Annual International ACM-SIGIR Conference on Research and Development in Information Retrieval*, pp. 129–136.
- Jing, H. and McKeown, K. (2000). “Cut and Paste Based Text Summarization.” In *Proc. of the 1st Meeting of the North American Chapter of the Association for Computational Linguistics*, pp. 178–185.
- Kan, M., McKeown, K., and Klavans, J. (2001). “Applying Natural Language Generation to Indicative Summarization.” In *Proc. of the 8th European Workshop on Natural Language Generation*, pp. 92–100.
- 笠原力弥, 山下洋一 (2001). “講演音声における重要文と韻律的特徴の関係.” 情報処理学会音声言語情報処理研究会報告, pp. 25–30. 35-5.
- 片岡明, 増山繁, 山本和英 (1999). “要約のための連体修飾節の“A の B”への言い換え.” 情報処理学会自然言語処理研究会報告, pp. 37–44. 133-6.
- 加藤直人, 浦谷則好 (2000). “放送ニュースを対象にした重要文抽出.” 言語処理学会第 6 回年次大会発表論文集, pp. 237–240.
- 川原裕美 (1989). “要約文のパラフレーズの諸相.” 佐久間まゆみ (編), 文章構造と要約文の諸相, pp. 141–167. くろしお出版.
- Knight, K. and Marcu, D. (2000). “Statistics-Based Summarization – Step One: Sentence Compression.” In *Proc. of the 17th National Conference on Artificial Intelligence*, p-p. 703–710.
- 小堀誠, 田村直良 (2000). “段落中の接続関係と段落間の重要度配分による文章要約.” 情報処理学会自然言語処理研究会報告, pp. 79–86. 136-11.
- Luhn, H. (1958). “The automatic creation of literature abstracts.” *IBM Journal of Research and Development*, 2 (2), pp. 159–165.
- 前田潤治, 高木啓伸, 福田健太郎, 浅川智恵子 (2001). “アノテーションに基づくウェブページのダイジェスト手法.” 電子情報通信学会福祉情報工学研究会報告, pp. 25–30. WIT-2001-14.

- Mani, I. (2001). *Automatic Summarization*. John Benjamins Publishing Company.
- Mani, I., Gates, B., and Bloedorn, E. (1999). "Improving Summaries by Revising Them." In *Proc. of the 37th Annual Meeting of the Association for Computational Linguistics*, pp. 558–565.
- Mani, I. and Maybury, M. (Eds.) (1999). *Advances in automatic text summarization*. MIT Press.
- Marcu, D. (1999). "The automatic construction of large-scale corpora for summarization research." In *Proc. of the 22nd Annual International ACM-SIGIR Conference on Research and Development in Information Retrieval*, pp. 137–144.
- Mathis, B., Rush, J., and Young, C. (1973). "Improvement of Automatic Abstracts by the Use of Structural Analysis." *Journal of the American Society for Information Science*, **24** (2), pp. 101–109.
- McKeown, K., Klavans, J., Hatzivassiloglou, V., Barzilay, R., and Eskin, E. (1999). "Towards Multidocument Summarization by Reformulation: Progress and Prospects." In *Proc. of the 14th National Conference on Artificial Intelligence*, pp. 453–460.
- 望主雅子, 荻野紫穂, 太田公子, 井佐原均 (2000). "重要文と要約の差異に基づく要約手法の調査." 情報処理学会自然言語処理研究会報告, pp. 95–102. 135-13.
- 難波英嗣, 奥村学 (1999). "書き換えによる抄録の読みやすさの向上." 情報処理学会自然言語処理研究会報告, pp. 53–60. 133-8.
- 奥村学, 難波英嗣 (1999). "テキスト自動要約に関する研究動向." 自然言語処理, **6** (6), pp. 1–26.
- 大塚敬義, 内海彰, 奥村学 (2001). "要約文生成における照応処理." 電子情報通信学会思考と言語研究会報告, pp. 19–26. TL-2001-4.
- Radev, D., Jing, H., and Budzikowska, M. (2000). "Centroid-based summarization of multiple documents: sentence extraction, utility-based evaluation, and user studies." In *Proc. of the ANLP/NAACL2000 Workshop on Automatic Summarization*, pp. 21–30.
- Reithinger, N., Kipp, M., Engel, R., and Alexandersson, J. (2000). "Summarizing Multilingual Spoken Negotiation Dialogues." In *Proc. of the 38th Annual Meeting of the Association for Computational Linguistics*.
- Stein, G., Strzalkowski, T., and Wise, G. (1999). "Summarizing Multiple Documents Using Text Extraction and Interactive Clustering." In *Proc. of the PACLING'99*, pp. 200–208.
- Tipster (1999). *Proceedings of The Tipster Text Program Phase III*. Morgan Kaufmann.
- 遠山義洋, 西田豊明 (2000). "話題構造の抽出と変形による対話録の自動要約." 2000年度人工知能学会全国大会論文集, pp. 157–160. 07-06.
- 上田良寛, 小山剛弘 (2000). "共通意味断片の抽出による複数文書要約." 言語処理学会第6回

年次大会発表論文集, pp. 360-363.

Witbrock, M. and Mittal, V. (1999). "Ultra-Summarization: A Statistical Approach to Generating Highly Condensed Non-Extractive Summaries." In *Proc. of the 22nd Annual International ACM-SIGIR Conference on Research and Development in Information Retrieval*, pp. 315-316.

Zechner, K. and Lavie, A. (2001). "Increasing the Coherence of Spoken Dialogue Summaries by Cross-Speaker Information Linking." In *Proc. of the NAACL2001 Workshop on Automatic Summarization*, pp. 22-31.

## 略歴

**奥村 学:** 1962年生. 1984年東京工業大学工学部情報工学科 卒業. 1989年同大学院 博士課程修了. 同年, 東京工業大学工学部情報工学科 助手. 1992年北陸先端科学技術大学院大学情報科学研究科 助教授, 2000年東京工業大学精密工学研究所 助教授, 現在に至る. 工学博士. 自然言語処理, 知的情報提示技術, 語学学習支援, テキストマイニングに関する研究に従事. 情報処理学会, 人工知能学会, AAI, 言語処理学会, ACL, 認知科学会, 計量国語学会 各会員. e-mail: oku@pi.titech.ac.jp, <http://oku-gw.pi.titech.ac.jp/~oku/>.

**難波 英嗣:** 1996年 東京理科大学理工学部電気工学科 卒業. 1998年 北陸先端科学技術大学院大学情報科学研究科 博士前期課程修了. 2001年 北陸先端科学技術大学院大学情報科学研究科 博士後期課程修了. 同年, 日本学術振興会 特別研究員. 2002年 東京工業大学精密工学研究所 助手. 現在に至る. 博士(情報科学). 自然言語処理, 特にテキスト自動要約の研究に従事. 情報処理学会, 人工知能学会, ACL, ACM 各会員. nanba@pi.titech.ac.jp, <http://oku-gw.pi.titech.ac.jp/~nanba/>

(2001年9月28日 受付)

(2001年12月21日 再受付)

(2002年4月5日 採録)