

# 研究終了報告書

## 「談話構造に基づく教師なし生成型要約」

研究期間：2019年10月～2022年3月

研究者：磯沼 大

### 1. 研究のねらい

近年、Webの発達と電子化された情報の飛躍的な増大により、自動文書要約のニーズが高まっている。例えばECサイト上に投稿された大量の商品レビューを要約することで、消費者の選択や製品企画に有用な情報を与えられる。自動要約は大量の文書からの知識抽出を容易にし、人の情報収集・意思決定の効率化・質的向上に貢献する技術である。

自動要約のアプローチは、要約に相応しい文や節を抽出する抽出型要約と、単語や句の言い換え・一般化を行う生成型要約に分けられる。生成型要約はより人手に近い自動要約を実現でき、その確立は自動要約研究の大きな目標である。一方で、生成型要約は見本となる要約(参照要約)を大量に要し、現実の文書の多くは参照要約の数が少なく、それらの用意に多大な労力を要することから、実用上の大きな障害となっている。

そこで、本研究は教師なし生成型要約手法を開発することによって、見本の要約が不要な汎用文書要約技術の実現に挑む。教師なしアプローチでは、要約の潜在表現を参照要約なしにいかに関得するかが鍵となる。本研究では文書の潜在構造、特にトピック構造に着目し、各トピックに関する要約文(トピック文)の潜在表現を得ることで、複数文で構成された要約を教師なしに生成する手法を開発する。例えば図1に示したあるレビューの要約は、food、place、serviceといった各観点について、様々な粒度の評価を述べている。トピック木構造を捉えた上でトピック文を生成することで、多様な粒度のトピックで構成された要約を生成できると考えた。

教師なし生成型要約は萌芽的な研究であり、本研究課題の申請時点にて申請者の研究を含め2例のみが報告されている。そのうち、本研究は文書に潜在するトピック構造に着目することで要約を生成する初の試みである。生成された要約と得られたトピック木構造との評価により、文書に潜在するトピック構造を捉えることが教師なし要約生成に有用であることを明らかにする。

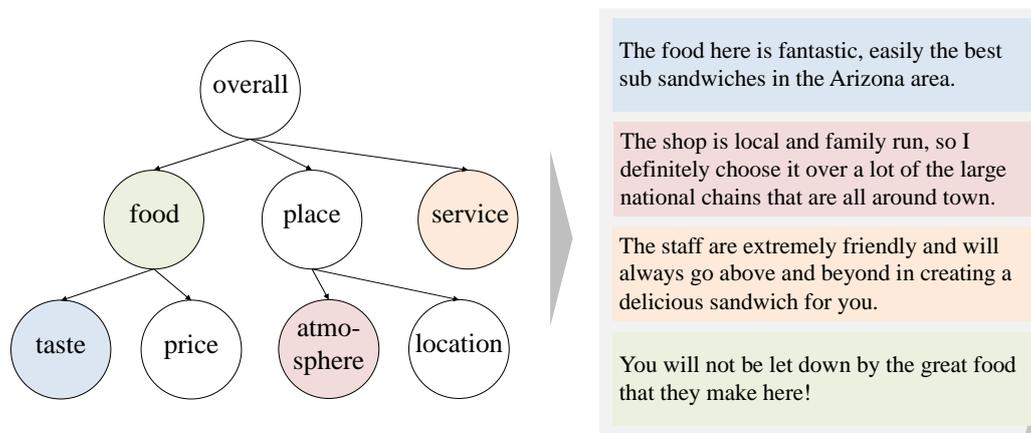


図1 あるレストランレビューの要約例と対応するトピック木構造

## 2. 研究成果

### (1) 概要

本研究課題では以下の過程を経て、トピック構造に基づく教師なし生成型要約手法の開発を行った。

前半では大規模文書に適用可能な木構造トピックモデルの開発に取り組んだ。本研究課題では木構造トピックモデルを用いて要約を生成するが、要約生成の学習には大量の文書が必要な一方、既存の木構造トピックモデルは大規模な文書に適用が困難である。そこで複数の文書を並列に学習できる木構造ニューラルトピックモデルを開発することで学習時間を短縮し、要約生成など大量の文書を要するタスクでの利用を可能にした。提案法は従来法と同等の解釈性を持つトピックとその木構造を得ながら、学習時間を約 15 倍短縮し、要約生成などのニューラルモデルとの一体的な学習を可能にするなど、木構造トピックモデルの応用可能性を広げた。本研究を取り纏めた論文は計算言語学分野のトップ国際会議 ACL2020 に採択された。

後半では、トピック構造に基づく教師なし要約生成手法の開発に取り組んだ。前半で開発した木構造トピックモデルにより、木構造上の各トピックに関する要約文の潜在表現を獲得し、要約を生成する手法を開発した。得られた要約と人手で作成した要約を比較評価した結果、既存の教師なし要約手法より元文書の内容を網羅し、かつ一貫性の高い要約が得られることが確認された。また、文の詳細度合いはその潜在分布の分散の大きさに依存し、潜在分布の分散が大きいほどより一般的な文が生成されるという、要約のみならず文生成タスク全体に有用な知見が得られた。以上の成果を取り纏めた論文は、計算言語学のトップジャーナル TAACL に採択された他、言語処理学会第 27 回年次大会で若手奨励賞を、情報処理学会第 246 回自然言語処理研究会で優秀研究賞および山下記念研究賞を受賞した。

### (2) 詳細

#### ・ 前半: 大規模文書に適用可能な木構造トピックモデルの開発

本研究では木構造トピックモデルを用いて要約を生成するが、要約生成の学習には大量の文書が必要な一方、既存の木構造トピックモデルは大規模な文書に適用が困難である。そこで前半では複数の文書を並列に学習できる木構造トピックモデルを開発することで学習時間を短縮し、要約生成など大量の文書を要するタスクでの利用を可能にした。

従来の木構造トピックモデルにて事後分布推定に用いられている崩壊型ギブスサンプリングや平均場近似は、学習性能や並列化の困難さから、大規模な文書への適用が困難である。また、要約生成とトピックの事後分布を一体的に学習することで、要約として有用なトピックとその木構造が得られるとより望ましい。

そこで本研究では、文書からトピック分布への写像をニューラルネットワークにより構成した木構造ニューラルトピックモデルを提案した。提案法は variational autoencoder (VAE) の枠組みによる学習が可能であり、前述の問題を解決可能である。

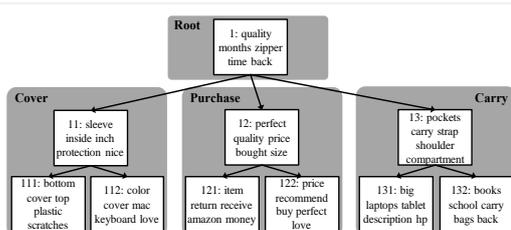


図 2 得られたトピック木構造

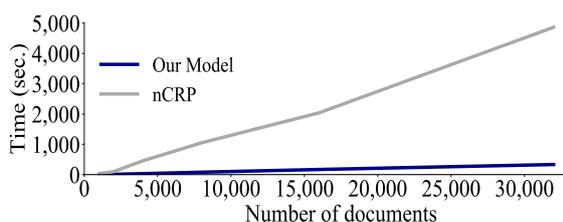


図 3 文書数(横軸)に対する学習時間(縦軸)

既存研究は LDA などのフラットなトピックモデルに VAE を適用しているが、無限木上のトピック分布への関数を有限のパラメータでどう構築するかは自明でない。そこで本研究は、親子間と兄弟間それぞれに再帰的な構造を持つ doubly-recurrent neural networks (Alvarez-Melis and Jaakkola, 2017) を用いて木構造上の stick-breaking process を表現することで、文書から無限木上のトピック分布への関数を構築した。

評価実験にて、提案法は既存の木構造トピックモデル(nCRP; Blei et al., 2010)とほぼ同等の一貫性を持つトピックと木構造を得た(図 2)。一方、学習時間は約 15 倍短縮され、大規模な文書に適用できることが示されたほか(図 3)、要約生成といった下流タスクのニューラルモデルとの一体的な学習が可能になるなど、木構造トピックモデルの応用可能性を広げた。

研究成果を国内会議にて発表したほか[3]、本研究を取り纏めた論文は計算言語学分野の国際会議 ACL2020 に採択された[1]。

#### ・ 後半:トピック構造に基づく教師なし要約生成手法の開発

本研究では、開発した木構造トピックモデルにより文書のトピック木構造を推定し、各トピックの要約文を生成する手法を開発した。トピック木構造では、文書から推定したトピックが木構造に配置され、子が親のサブトピックとなる構造を持つ。それらのトピックから要約として相応しい詳細度合いのトピックを選択し、各トピックに関する要約文を生成することで、意見文書の要約が教師なしに得られることを示した。

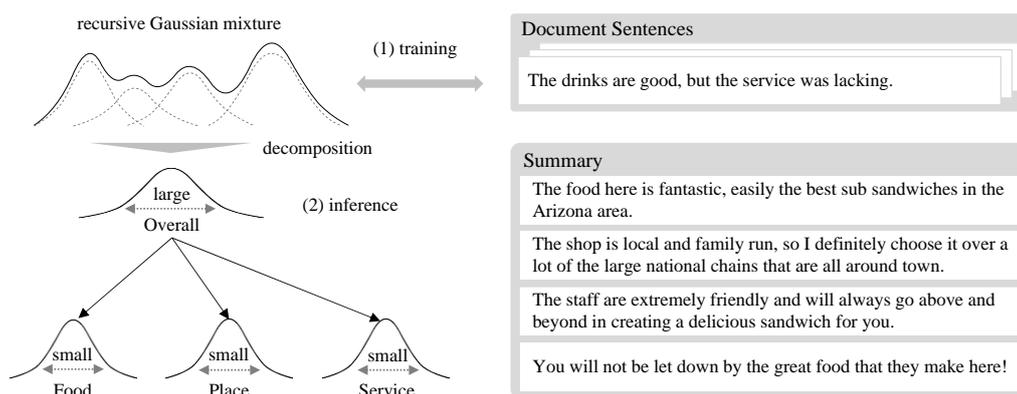


図 4 提案法の概要図。(1) 学習フェーズでは、文書の文の潜在変数が混合ガウス分布に従うと仮定し、そのパラメータを学習する。(2) 要約生成フェーズでは、混合ガウス分布を構成する各単峰ガウス分布からトピック文を生成し、その集合を要約として出力する。潜在分布の分散は葉に近づくほど小さくなるようモデル化されており、葉に近づくにつれより具体的なトピックに関する要約文が出力される。

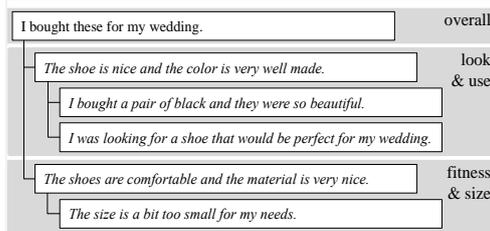


図 6 冠婚葬祭用シューズのレビューから得られた要約(トピック文の集合)の例

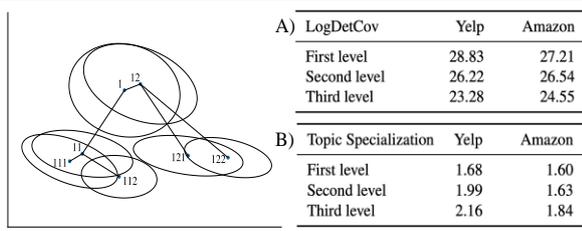


図 5 トピック文の潜在分布の可視化(PCA)。葉に近づくにつれ分布の分散が小さくなる一方(表 A)、トピック文の内容は詳細になる(表 B)

トピック文生成の文脈では、文書中の文の潜在分布を混合ガウス分布で表現することで、その構成要素である各単峰ガウス分布がトピック文の潜在分布として機能することを明らかにした既存研究が存在する(Wang et al., 2019)。一方、本研究のように多様な詳細度合いを持つトピック文を生成するためには、文の詳細度合いを潜在空間上でモデル化する必要がある、その方法は明らかでない。

そこで本研究では、トピック文の潜在表現をガウス分布で表現する際に、子の分布の分散が親よりも小さくなるようにモデルを構築することで、根からは抽象的な文を、葉に近づくにつれより詳細な文を生成することを試みた。単語の潜在分布としてガウス分布を用いる Gaussian word embedding では、「犬」のような具体的な単語は、「動物」といった一般的な単語よりも小さい分散を持つことが示されている(Vilnis et al., 2015)。文においても同様に、具体的な文は意味の分散が小さいため、その潜在分布は分散が小さいと考えられる。

子の分布の分散を親よりも小さくするために、本研究では再帰的混合ガウス分布(再帰的 GMM)を文書中の文の潜在表現の事前分布として導入した。再帰的 GMM は、木構造上の各トピックに対応するガウス分布で構成され、子の事前分布に親の事後分布を設定して構築される。これにより、根の潜在分布の分散は大きく、葉に近づくにつれ分散は小さくなる。

評価実験では、提案法の要約性能は最新の教師なし生成型要約手法(Bražinskis et al., 2020)と競合することを確認した。図 5 に示すように、提案法は“look & use”や“fitness & size”といった多様なトピックを捉えており、人手評価では要約の informativeness や元の文書の被覆率といった観点で既存手法を上回ることを確認した。また、トピック文の詳細度合いはその潜在分布の分散の大きさに依存し、根の文の潜在分布は分散が大きく一般的な文が生成される一方、葉に近づくにつれ分散が小さくなり具体的な文が生成されるといった特性を確認した(図 6)。これは単語の潜在表現にガウス分布を用いた Gaussian word embedding にて報告された特性と類似しており、要約のみならず、質問応答や対話生成などの文の詳細度合いを考慮する他タスクにも有用な知見である。

以上の成果を取り纏めた論文は、計算言語学のトップジャーナル TAACL に採択された他[2]、言語処理学会第 27 回年次大会で若手奨励賞を[5][7]、情報処理学会第 246 回自然言語処理研究会で優秀研究賞および山下記念研究賞を受賞した[4][6][8]。

### 3. 今後の展開

これまで 2 年間の研究計画では、トピック構造に基づく教師なし文書要約技術の確立を目指しており、その目標は達成されたものと認識している。

しかし、教師なし文書要約の最終的な目標は汎用的な自動文書要約の実現であり、これまで評価実験で用いてきた商品レビュー以外にも適用できるのかという工学的な関心が残る。特にこれまで自動文書要約技術が適用されてこなかった領域を、トピック構造を捉える本アプローチによって開拓することを目指したい。

一方、計算言語学的な観点からみると、トピック構造を捉えることが要約の質にどのように寄与するのか、実験的には明らかにしたものの理論的には明らかにできていない。要約に求められる性質(非冗長性、原文書との関連性、情報量の多さ)を定式化した上で、それらとトピック構造を捉えることの数理的なつながりを議論したい。

### 4. 自己評価

- 研究目的の達成状況

2 年間の研究計画では、トピック構造に基づく教師なし文書要約技術の確立を目指しており、その目標は達成されたと認識している。

- 研究の進め方(研究実施体制及び研究費執行状況)

研究費は研究補助を担う学生の謝金や雇用経費に多く充てられており、リヴァプール大学ダヌシカ・ボレガラ教授や当該学生など、研究室内外の研究者との共同研究によって、研究が大きく加速した。

- 研究成果の科学技術及び社会・経済への波及効果

本研究で開発した教師なし要約生成は、自動要約技術の実用化に重要な技術であり、社会・経済的にも重要な意義を持つと認識している。これまで自動文書要約技術が適用されてこなかった領域を今後開拓することによって、研究成果の社会的意義をより直接的に発信していきたい。

- その他

ACT-X 開始以前は研究者としての個を確立していくという意識が欠けていたが、領域内の研究者やアドバイザーからの触発を通じて、いかに個を確立していくか常に自問するようになったことは、ACT-X で通じて得られた最も大きな価値の一つであった。

## 5. 主な研究成果リスト

### (1) 代表的な論文(原著論文)発表

研究期間累積件数:2件

[1] Masaru Isonuma, Junichiro Mori, Danushka Bollegala, and Ichiro Sakata. Tree-Structured Neural Topic Model. In Proceedings of the 58th Annual Meeting of the Association for Computational Linguistics (ACL), pp.800-806, 2020.

概要 本研究では複数の文書を並列に学習できる木構造ニューラルトピックモデルを開発することで学習時間を短縮し、大量の文書に対するトピック構造の推定を可能にした。提案法は従来法と同等の解釈性を持つトピックとその木構造を得ながら、学習時間を約15倍短縮し、ニューラルモデルを用いた要約生成など下流タスクとの一体的な学習を可能にするなど、木構造トピックモデルの応用可能性を広げた。

[2] Masaru Isonuma, Junichiro Mori, Danushka Bollegala, and Ichiro Sakata. Unsupervised Abstractive Opinion Summarization by Generating Sentences with Tree-Structured Topic Guidance. Transactions of the Association for Computational Linguistics (TACL), vol.9, pp.945-961, 2021.

概要 本研究では開発した木構造トピックモデルにより文書のトピック木構造を推定し、各トピックの要約文を生成する手法を開発した。トピック木構造では、文書から推定したトピックが木構造に配置され、子が親のサブトピックとなる構造を持つ。それらのトピックから要約として相応しい粒度のトピックを選択し、各トピックに関する要約文を生成することで、意見文書の要約が教師なしに得られることを示した。

### (2) 特許出願

研究期間全出願件数:0件(特許公開前のもも含む)

### (3) その他の成果(主要な学会発表、受賞、著作物、プレスリリース等)

#### 国内会議発表

[3] 磯沼 大, 森 純一郎, ボレガラ ダヌシカ, 坂田 一郎. 木構造ニューラルトピックモデル. 言語処理学会第26回年次大会, 2020.

[4] 磯沼 大, 森 純一郎, ボレガラ ダヌシカ, 坂田 一郎. 潜在的なトピック構造を捉えた生成型教師なし意見要約. 情報処理学会第246回自然言語処理研究会, 2020.

[5] 磯沼 大, 森 純一郎, ボレガラ ダヌシカ, 坂田 一郎. トピック文生成による教師なし意見要約. 言語処理学会第27回年次大会, 2021.

#### 受賞

[6] 優秀研究賞, 情報処理学会第246回自然言語処理研究会, 2020.

[7] 若手奨励賞, 言語処理学会第27回年次大会, 2021.

[8] 山下記念研究賞, 情報処理学会, 2021.