文の分散表現に基づく小説の自動分割とストーリー展開の解析

Novel Segmentation Method based on the Distributed Representation of Sentences and Analysis Method of Story Developments

> 福田 清人^{*1} 森 直樹^{*1} 岡田 真^{*1} Kiyohito Fukuda Naoki Mori Makoto Okada

> > *1大阪府立大学 工学研究科

Graduate School of Engineering, Osaka Prefecture University

Recently, the attempts to reproduce the mechanisms of human intellectual activities have attracted interest in artificial intelligence fields. The narrative creation is one of them. It is necessary for narrative creation and creative support by the computer to make it to understand human creations and their stories. However, there are few studies on story analysis by the computer. In this study, we propose the segmentation method of novels based on the distributed representation of sentences and the analysis method of story developments. As a result of computational experiments, we confirmed that the effectiveness of the proposed methods.

1. はじめに

近年,人の知的活動の仕組みを計算機上で再現する試みが人 工知能の分野で広く行われ,大きな注目を集めている.人の知 的活動の1つに物語の創作がある.人の感性に基づく創作物 である物語はストーリーと表現媒体の2つの要素から構成さ れ,それらの組合せによって小説や漫画などに枝分かれしてい く.ここで,ストーリーは物語の内容であり,表現媒体はその 内容をどのような形で表現するかという表現方法である.計算 機による物語の自動生成の研究 [Pěrez 01][Ueno 14] は数多く 報告されている.また,近年では人と機械による共同作業に関 する研究も注目を集めており,人と計算機による物語の共同創 作の研究 [上原 11] や創作支援の研究 [葛井 17] も報告されて いる.どちらの研究においても,人の創作物を計算機に理解さ せることが非常に重要となる.

物語の自動生成や創作支援を実現するためには,既存の物語 を解析し,人が物語を創作するうえで必要な知識や技術を計算 機が理解可能な形で獲得する必要がある.具体的には,機械学 習技術に基づく既存の物語に対する解析による有用な情報の抽 出は必要不可欠な技術である.しかしながら,物語の解析に関 する研究は,専門家の経験則に基づいて人手で情報を抽出する 研究[佐藤 10]が主流であり,計算機による数値的な情報抽出 や解析に関する研究はほとんど報告されていない.

なお,本研究ではストーリーに焦点を当てる.表現媒体が時 代とともに姿を大きく変えることがある一方,古くから存在す るストーリーの典型的な構造が現代でもしばしば使用されるこ とから,表現媒体と比べてストーリーの方が時間経過に対して ロバストなためである.また,本研究における解析対象には小 説を用いる.

以上を背景として、本研究ではストーリーを計算機に理解させるための第一歩として、文意を考慮した小説の自動分割手法 およびストーリー展開の解析手法を提案する.文意を考慮した 文の分散表現に基づき、小説文をストーリーが展開する部分で 自動分割する.また、自動分割された複数の小説文をシーンと みなし、各シーンのベクトルを用いてストーリー展開が類似し た部分を発見する.

2. 関連研究

本研究と関連のあるいくつかの研究について説明する.

2.1 テキストセグメンテーションに関する研究

テキストデータをトピックなどの意味的なまとまりに分割 するテキストセグメンテーションに関する代表的な手法に TextTiling[Hearst 97] がある. TextTiling はテキスト中の ある2文間を基準として、その前後の文をあらかじめ設定した 窓幅の分だけそれぞれ取得し、得られた前後の文章に対して単 語の出現頻度ベクトルの類似度を計算する. この操作を基準と なる2文を動かしながら実行し、得られた類似度の変化から文 章境界を推定する手法である. TextTiling は文章内に出現す る単語に基づいてセグメンテーションするため、短い文章を対 象とした場合には有効に機能しないことが知られている.

2.2 物語の解析に関する研究

物語の解析に関する研究では, 星新一の作品を構造分析の考 えに基づき, テキストの時系列に着目して物語のパターン抽出 をする研究が報告されている [佐藤 10]. しかしながら, 物語の パターンを抽出するためにはテキストを抽象化して分類する必 要があるため, 人手によってしか解析できないという問題点が 存在する.

3. 提案手法

本研究では,物語の中でもストーリーという要素に着目して 小説を解析する.ここで,ストーリーをイベントや登場人物の 行動,場所移動に伴う物語中の一連の状態遷移の時系列である と定義する.小説をいくつかの文章の集合であると仮定する と,ある連続した2つの文章間の差が状態遷移であり,冒頭か ら末尾までの連続した2文章間の差の分布がストーリーであ るといえる.そこで,小説中の文章を分散表現化して文章ベク トルを得た場合,小説は文章ベクトルの時系列集合とみなすこ とができる.また連続する2つの文章ベクトルに何らかの演 算子を適用した結果がその2文章間での差であり,物語内の状 態遷移を表しているといえる.そのため,小説におけるストー リーは冒頭から末尾までの連続した2つの文章ベクトルにあ る演算子を適用した結果の時系列集合であると定義できる.

以上の観点から本研究では,文の分散表現に基づく小説文の 自動分割手法およびストーリー展開の解析手法を提案する.な

連絡先: 福田清人, 大阪府立大学 工学研究科, 大阪府堺市中区学 園町 1-1, 072-254-9273, fukuda@ss.cs.osakafu-u.ac.jp

お, 文の分散表現の獲得には, これまでに提案してきた文の分 散表現の獲得手法 [Fukuda] を改良した手法を用いる.

3.1 文の分散表現を用いた小説文の自動分割

TextTiling の考え方を基にして,小説文の各文の分散表現 に対して類似度を計算し,類似度が極大となる2文を結合して いく操作を,セグメント数が任意の数となるまで繰り返すこと で小説を自動分割する手法を提案する.ここで1文単位での 類似度計算をすると,機械的な文分割により分割されてしまっ た不適切な文の前後を分割点と推測してしまう可能性がある. そこである1文に対して,その1文と前後窓幅分を含む文の 分散表現の平均を類似度計算に用いるベクトルとするスムー ジング手法を導入する.図1および図2に小説文のセグメン テーション手法の概要およびスムージング手法を示す.以下に 文の分散表現を用いた小説文のセグメンテーション手法のアル ゴリズムを示す.

- 1. 獲得したいセグメント数を $N_{\rm s}$,スムージング幅を $N_{\rm w}$ と する.ここで、本節で用いるスムージング手法では基準と なる文に対してその前後の $N_{\rm w}$ 文を含む $2N_{\rm w}$ + 1 文を まとめてスムージングする.
- 2. 解析する小説を M 文の文集合とする.
- 3. 小説中の各文に対して, 文の分散表現の獲得手法により文 の分散表現 *s*_i(*i* = 1,2,...,*M*) を獲得する.
- 4. 各セグメントに対応したセグメントベクトルを d_j ($j = N_w + 1, N_w + 2, \cdots, M N_w$), セグメントベクトルの 集合を $\mathcal{D} = \{d_{N_w+1}, d_{N_w+2}, \cdots, d_{M-N_w}\}$ とする.ま た,各セグメントに含まれる文数を b_j ,この文数の集合を $\mathcal{B} = \{b_{N_w+1}, b_{N_w+2}, \cdots, b_{M-N_w}\}$ とする.ここで,

$$\boldsymbol{d}_{j} = \frac{1}{2N_{\mathrm{w}}+1} \sum_{k=j-N_{\mathrm{w}}}^{j+N_{\mathrm{w}}} \boldsymbol{s}_{k}$$
(1)

$$b_j = \begin{cases} N_{\rm w} + 1 & (j = N_{\rm w} + 1, M - N_{\rm w}) \\ 1 & (\text{otherwise}) \end{cases}$$
(2)

である.

5. セグメントベクトル集合の連続した 2 つのセグメントベ クトル d および d' の類似度 $f_{sim}(d, d')$ を以下の式に 従って計算する. ここで, γ は減衰率であり, セグメン トが長文になりすぎないよう制御するための可調整パラ メータである. γ は 0 < γ < 1 を満たす実数である. ま た, b および b' はそれぞれセグメントベクトルに対応し たセグメントに含まれる文数である.

$$f_{\rm sim}\left(\boldsymbol{d}, \boldsymbol{d}'\right) = \gamma^{b+b'-2} \left(1 + \frac{\boldsymbol{d} \cdot \boldsymbol{d}'}{|\boldsymbol{d}||\boldsymbol{d}'|}\right) \tag{3}$$

- 6.5 で求めた類似度が最大となった2つのセグメントベクトルを $d_{\rm m}$ および $d_{\rm m'}$ とし、それぞれに対応するセグメントに含まれる文数をそれぞれ $b_{\rm m}, b_{\rm m'}$ とする.
- *d*_m および *d*_{m'} に対応するセグメントを結合し、1 つの セグメントとする. その後、以下の操作を適用することで 各値を更新する. ここで、記号 '→' は 左式を右式で更新



図 1: 小説文のセグメンテーション手法の概要



図 2: 小説文のスムージング手法の概要

する操作を表す.

$$d_m = \frac{1}{2N_{\rm w} + b_{\rm m} + b_{\rm m'}} \sum_{k=m-N}^{m'+b_{\rm m'}+N_{\rm w}-1} s_k \quad (4)$$

$$b_{\rm m} \rightarrow b_{\rm m} + b_{\rm m'}$$
 (5)

$$\mathcal{B} \to \mathcal{B} \setminus \{b_{\mathbf{m}'}\} \tag{6}$$

$$\mathcal{D} \rightarrow \mathcal{D} \setminus \{ \boldsymbol{d}_{\mathrm{m}'} \}$$
 (7)

- 8. $|\mathcal{D}| > N_{\rm s}$ ならば, 5 に戻る.
- |D| = N_sの時, D および D の各要素に対応したセグメントを獲得する.

3.2 文の分散表現を用いたストーリー展開の解析

3.1 節で説明した自動分割手法により得られたセグメントと そのセグメントに対応したセグメントベクトルを用いてストー リー展開を解析する.本節では,連続したセグメントに対応し たベクトル間の差分を計算し,得られた差分ベクトルを1つの ストーリー展開とみなすことでストーリーを多次元数値空間上 で表現する.以下にストーリー展開の解析手順を示す.

- 1. 解析対象とする作品を複数用意する.
- 3.1 節の自動分割手法により作品を自動分割し、各作品の セグメントとそれに対応したセグメントベクトルを獲得 する.
- 3. 各作品ごとに, 連続した 2 つのセグメントベクトルの差 分を計算する.

表 1: 実験 1 における実験条件

$m_{\rm max}$	80
	太宰治 「走れメロス」
	太宰治 「黄金風景」
使用作品	芥川龍之介 「蜘蛛の糸」
	芥川龍之介 「藪の中」
	エドガー・アラン・ポー 「黒猫」

- 得られた差分ベクトルに対して、他作品から得られた差分 ベクトルとのコサイン類似度を計算する.
- 得られたコサイン類似度や差分ベクトルの各要素の割合 に基づいて,作品間でのストーリーの類似性や,差分ベク トルとストーリー展開の関係性などを可視化しつつ解析 する.

4. 実験

提案手法の有効性を確認するため、いくつかの実験をした. また,提案手法を用いていくつかのストーリーを実際に解析す ることで,得られる情報の特徴や傾向についての知見を得る. 実験1~3まで実施したが,紙面の関係上,実験1のみ示す. 実験2および実験3については発表時に述べる.

4.1 実験 1

文の意味的な類似性を考慮した文の分散表現を用いて,小説 のストーリー展開を可視化することができるかを確認する.

小説の各文に対して文の分散表現を獲得し, 次元圧縮手法で ある t-distributed Stochastic Neighbor Embedding (t-SNE) を用いて 2 次元空間に写像することで可視化する.

4.2 実験1の実験条件

表1に実験1の実験条件を示す.実験1では青空文庫から 取得した各作品を文単位に分割し,その中から m_{max} 単語以下 の文を用いた.

4.3 実験1の結果と考察

図3~図7に各作品のt-SNEによる可視化の結果を示す. ここで,各図において独立にt-SNEを適用しているため,それ ぞれの図の軸に関係性はないことに注意する.また,t-SNEは データ間の距離を確率分布で表現することで次元を圧縮する ため,次元圧縮後の各軸には意味がないことにも注意を要を要 する.

図3を見ると、各文の分散表現が2つの分布パターンに属 していることがわかる.これは、「走れメロス」という作品は 町での王様との会話シーンと、村に戻ってまた町に戻るという 移動シーンの2つに分けることができることから妥当である と考えられる.

図4および図5を見ると、すべての文の分散表現が次元圧 縮後の2次元空間において、ある直線上に分布していることが わかる.このような分布をとる原因として、「黄金風景」およ び「蜘蛛の糸」は場面転換時の風景描写などが淡白であること や、作品全体を通して使用される文体や人物の口調なども一定 であることが考えられる.

図6を見ると、「藪の中」の各文の分散表現は先頭から末尾 に向かって一定方向に展開していることがわかる.このことか ら、「藪の中」は物語の後半部分で前半部分について回想する こともなく、ストーリーが次々に展開していくと考えられる.



図 3: t-SNE による「走れメロス」の可視化結果



図 4: t-SNE による「黄金風景」の可視化結果

実際に,「藪の中」という作品はある男が殺される事件に対し て,目撃者や容疑者,被害者本人の霊などが事件について語る というストーリーであり,それぞれが自身の体験を語るだけで, 各人物が関わりあうシーンが存在しない.

図7を見ると,文の分散表現が一様に分布していることがわ かる.また,「黒猫」という作品は主人公の身の回りで起きる 事件と主人公の内面が交互に書かれている作品であり,類似し たストーリー展開を作品内で何度か繰り返す内容になってい る.これらのことから,似たストーリー展開を繰り返すことで, 時系列という観点から文の分散表現を可視化した結果,一様に 分布してしまっていると考えられる.

以上の結果から、小説のストーリーやその展開に関して、文 の分散表現を用いた多次元数値空間上で作品の種類や特徴を解 析可能であると考えられる.

5. むすび

本研究では,文の意味的な類似性を考慮した文の分散表現に 基づく小説の自動分割手法とストーリー展開の解析手法を提案 した.自動分割手法では TextTiling の考え方を基にして,文 の分散表現間の類似度が極大な部分を結合することで小説を シーンごとに自動分割した.ストーリー展開の解析手法では,



図 5: t-SNE による「蜘蛛の糸」の可視化結果



図 6: t-SNE による「藪の中」の可視化結果

自動分割手法で得られたセグメント間の差分をストーリー展開 とみなし,差分ベクトルの類似度からストーリー展開の類似性 を解析した.提案手法の有効性を確認するためのいくつかの実 験により,以下の知見が得られた.

- 小説文の分散表現を次元圧縮手法である t-SNE によって 可視化することで、小説のストーリーやその展開について 多次元数値空間上で解析することができる。
- 提案手法を用いることで、人手によるアノテートに頼ることなく小説文を意味を考慮したシーン単位に自動分割することができる。
- 提案手法により、ストーリー展開の類似性だけではなく、 文章構成の類似性も取得することができる。

今後の課題として提案手法の根幹となる文の分散表現の性 能向上は最重要課題である. 階層的 LSTM や Attention 機構 のような自然言語処理における有効性が示された技術を手法に 導入することで更なる性能向上が期待される. また, 自動分割 手法におけるスムージングの窓幅やシーンの分割数など解析対 象の作品ごとに異なる可調整パラメータを最適化するための手 法について検討する必要がある.



図 7: t-SNE による「黒猫」の可視化結果

謝辞

本研究は一部,日本学術振興会科学研究補助金基盤研究(C) (課題番号 26330282)の補助を得て行われたものである.

参考文献

- [Fukuda] Fukuda, K., Mori, N., and Matsumoto, K.: A Novel Sentence Vector Generation Method Based on Autoencoder and Bi-directional LSTM, in Prieta, de la F., Omatu, S., and Fernández-Caballero, A. eds., Distributed Computing and Artificial Intelligence, 15th International Conference, DCAI 2018, Toledo, Spain, 20-22 June 2018, Vol. 800 of Advances in Intelligent Systems and Computing, pp. 128–135
- [Hearst 97] Hearst, M. A.: TextTiling: Segmenting Text into Multi-paragraph Subtopic Passages, Comput. Linguist., Vol. 23, No. 1, pp. 33–64 (1997)
- [Pěrez 01] Pěrez, y R. P. and Sharples, M.: MEXICA: A computer model of a cognitive account of creative writing, Journal of Experimental & Theoretical Artificial Intelligence, Vol. 13, No. 2, pp. 119–139 (2001)
- [Ueno 14] Ueno, M., Mori, N., and Matsumoto, K.: 2-Scene Comic Creating System Based on the Distribution of Picture State Transition, pp. 459–467, Springer International Publishing, Cham (2014)
- [葛井 17] 葛井 健文, 上野 未貴, 井佐原 均: 質問集合とグラフ に基づく物語全体の流れを管理可能な創作支援システムの 提案, 第 31 回人工知能学会全国大会発表論文集 (2017)
- [佐藤 10] 佐藤 知恵, 村井 源, 徃住 彰文: 星新一ショートショー ト文学の物語パターン抽出, 情報知識学会誌, Vol. 20, No. 2, pp. 123–128 (2010)
- [上原 11] 上原 大輝, 出水 ちあき, 宮里 洸司, 神里 志穂子, 野口 健太郎: J-030 子どもの思考プロセス把握における物語自作システムの有効性検証 (HCS(2),J 分野:ヒューマンコミュニケーション&インタラクション), 情報科学技術フォーラム 講演論文集, Vol. 10, No. 3, pp. 597–600 (2011)