

平成29年度 卒業論文

分散表現を用いた 話題変化判定

提出日 平成29年9月19日

所属 名古屋工業大学 情報工学科

指導教員 伊藤 孝行 教授

入学年度 平成29年度入学

学籍番号 26115162

氏名 芳野 魁

論文要旨

近年、Web 上での大規模な議論活動が活発になっているが、現在一般的に使われている ”2ちゃんねる” や ”Twitter” といったシステムでは整理や収束を行うことが困難である。困難である原因として、議論の管理を行う者がいないことが挙げられる。議論を収束させるには議論のマネジメントを行う人物が必要である。大規模意見集約システム COLLAGREE ではファシリテーターと呼ばれる人物が議論のマネジメントを行っている。しかし、ファシリテーターは人間であり、長時間に渡って大人数での議論の動向をマネジメントし続けるのは困難である。COLLAGREE で大規模な議論を収束させるためには、ファシリテーターが必要な時にだけ画面を見るようにして画面に向き合う時間を減らす工夫があることが望ましい。ファシリテーターが画面を見るべきタイミングは議論の話題が変化したときである。以前の議論の内容から外れた発言がされた時、ファシリテーターが適切に発言することで、脱線や炎上を避けて議論を収束させることができる。すなわち、ファシリテーターの代わりに自動的に議論中の話題の変化を事前に判定することが求められている。現在、COLLAGREE 上で使用されている議論支援システムは投稿支援システムと議論可視化システムの2つに大別できる。投稿支援システムはポイント機能やファシリテーションフレーズ簡易投稿機能のように、ユーザーが投稿をする際に何らかの補助やリアクションを行う。現行の機能では選択肢の提示に留まっており、作業量を減らすことには繋がりにくい。一方、議論可視化システ

ムは議論ツリーやキーワード抽出のように、ユーザーにスレッドとは異なる議論の見方を提供する。現行の機能では議論を見やすくすることに重点が置かれており、議論の把握の助けにはなるが画面に向き合う時間を減らすことにはなりにくい。むしろ、作業量を増やすことになり得る機能もある。近年、自然言語処理の分野において分散表現が多くの研究で使われており、機械翻訳を始めとする単語の意味が重要となる分野で精度の向上が確認されている。分散表現を用いることで、人間に近い精度で話題の変化を観測することが可能となる。以上のような背景を踏まえて、分散表現を用いて、話題の変化を観測し、話題の変化が確認された時にファシリテーターに伝えることが望ましい。話題の変化の観測は、発言中に現れる単語の関連度合いの計算と見なすことができる。分散表現を用いることで単語間の類似度を求めることができる、値が大きいほど単語がそれぞれ類似した実数ベクトルであることを表す。単語 A と単語 B の実数ベクトルが類似しているとは、単語 A と共に使われることの多い単語と単語 B と共に使われることの多い単語が多く共通していることを示す。故に、分散表現を使って単語の関連度を計算することができる。発言文から単語を選ぶ際には自動要約を用いる。発言文から重要でない単語を取り除くことで関連度の計算の精度を高めることが可能となる。本論文では、分散表現を用いて議論中での発言に含まれる単語の関連度を計算し、話題の変化を観測する手法を提案する。提案手法は、既存の抽出的要約手法を用いて選ばれた単語の関連度を計算する手法、Seq2Seq による生成的要約を用いて生成された単語の関連度を計算する手法、オントロジーを用いて求められた単語の関連度を計算する手法の 3 つである。提案した 3 つの手法により、議論中の話題の変化の観測の評価実験を行い、各手法の評価を行う。評価実験によって、提案手法を用いることで人間の代わりに自動的に話題の変化を観測できるこ

とを確認する.

目 次

論文要旨	1
目次	4
図目次	8
表目次	9
第 1 章 序論	11
1.1 研究の背景	11
1.2 研究の目的	13
1.3 本論文の構成	13
第 2 章 関連研究	15
2.1 序言	15
2.2 COLLAGREE	16
2.2.1 概要	16
2.2.2 ファシリテーター	17
2.3 議論支援	17
2.4 話題遷移検出	18
2.4.1 テキストセグメンテーション	19

2.4.2	トピックモデル	21
2.5	重み付け	23
2.5.1	Okapi BM25	23
2.5.2	LexRank	25
2.6	分散表現	27
2.6.1	単語文脈行列	29
2.6.2	word2vec	30
2.6.3	fastText	31
2.6.4	分散表現を使用した話題関連研究	32
2.7	結言	34
第3章	通知システム	35
3.1	序言	35
3.2	システムの動作の流れ	35
3.3	システム詳細	37
3.3.1	発言データ	37
3.3.2	発言間の類似度計算	38
3.4	結言	40
第4章	発言内容の類似度計算	41
4.1	序言	41
4.2	前処理	41
4.2.1	MeCab	42
4.2.2	重み付け	43

4.3	類似度計算	44
4.3.1	単語抽出	44
4.3.2	分散表現による類似度計算	44
4.4	結言	45
第 5 章	評価実験	47
5.1	序言	47
5.2	対象データ	48
5.2.1	議論データ	48
5.2.2	評価データ	49
5.3	実験設定	52
5.3.1	パラメーター	52
5.3.2	比較手法	52
5.3.3	評価指標	53
5.4	実験結果	54
5.5	結言	54
第 6 章	結論	55
6.1	序言	55
6.2	今後の課題・展望	55
6.3	本研究のまとめ	55
	謝辞	57
	付 録 A 発表 (予定) 論文一覧	65

A.1 発表(予定)論文一覧	65
A.2 投稿論文一覧	65
付 録 B 市民共創知研究会 投稿論文	67

目 次

1.1	議論ツリー	12
2.1	COLLAGREE のスレッドの例	16
2.2	TextTiling の流れ	19
2.3	TextTiling のグラフ	20
2.4	ユニグラムモデル	22
2.5	混合ユニグラムモデル	22
2.6	LDA	22
2.7	類似度グラフの例	26
2.8	単語文脈行列	29
2.9	CBOW	30
2.10	skip-gram	30
2.11	話題展開システムの構造	32
2.12	TweetSift によるトピック予測のワークフロー	33
3.1	システムの流れ	36
4.1	解析結果	42

表 目 次

3.1	発言データ	38
5.1	パラメーターの設定	53
5.2	実験結果	54

第1章

序論

1.1 研究の背景

近年，Web 上での大規模な議論活動が活発になっているが，現在一般的に使われている ”2ちゃんねる” や ”Twitter” といったシステムでは整理や収束を行うことが困難である．困難である原因として，議論の管理を行う者がいないことが挙げられる．つまり，議論を整理・収束させるには議論のマネジメントを行う人物が必要である．大規模意見集約システム COLLAGREE[16] ではファシリテーターと呼ばれる人物が議論のマネジメントを行っている．しかし，ファシリテーターは人間であり，長時間に渡って大人数での議論の動向をマネジメントし続けるのは困難である．COLLAGREE で大規模な議論を収束させるためには，ファシリテーターが必要な時には画面を見るようにして，他の時は見なくても済むようにすることで画面に向き合う時間を減らす工夫があることが望ましい．ファシリテーターが画面を見るべきタイミングは議論の話題が変化したときである．以前の議論の内容から外れた発言がされた時，ファシリテーターが適切な発言をすることで，脱

線や炎上を避けて議論を収束させることができる。すなわち、ファシリテーターの代わりに自動的に議論中の話題の変化を観測することが求められている。現在、COLLAGREE 上で使用されている議論支援システムは「(1) 投稿支援システム」と「(2) 議論可視化システム」の2つに大別できる。投稿支援システムはポイント機能やファシリテーションフレーズ簡易投稿機能のように、ユーザーが投稿をする際に何らかの補助やリアクションを行う。現行の機能では選択肢の提示に留まっており、作業量を減らすことには繋がりにくい。一方、議論可視化システムは議論ツリーやキーワード抽出のように、ユーザーにスレッドとは異なる議論の見方を提供する。1.1 に議論ツリーの例を示す。現行の機能では議論を見やすくすること

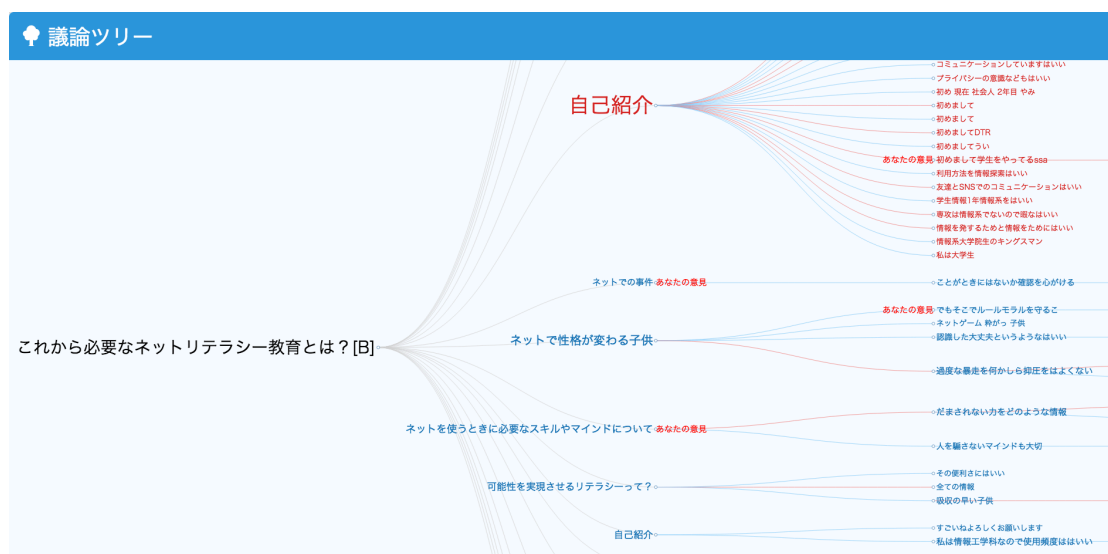


図 1.1: 議論ツリー

に重点が置かれており、議論の把握の助けにはなるが画面に向き合う時間を減らすことにはなりにくい。むしろ、作業量を増やすことになり得ることもある。従って、現行の支援機能ではファシリテーターの作業量の減少には繋がりにくい。近年、自然言語処理の分野において分散表現が多くの研究で使われており、機械翻

訳を始めとする単語の意味が重要となる分野で精度の向上が確認されている。分散表現を用いることで、人間に近い精度で話題の変化を観測することが可能となる。以上のような背景を踏まえて、分散表現を用いて、話題の変化を観測し、話題の変化が確認された時にファシリテーターに伝えることが望ましい。話題の変化の観測は、発言中に現れる単語の類似度の計算と見なすことができる。分散表現を用いることで単語間の類似度を求めることができる、値が大きいほど単語がそれぞれ類似した実数ベクトルであることを表す。単語 A と単語 B の実数ベクトルが類似しているとは、単語 A と共に使われることの多い単語と単語 B と共に使われることの多い単語が多く共通していることを示す。故に、分散表現を使って単語の類似度を計算することができる。発言文から単語を選ぶ際には自動要約を用いる。発言文から重要でない単語を取り除くことで関連度の計算の精度を高めることが可能となる。要約の手法としては okapi BM25 [5] と LexRank を組み合わせた抽出的要約手法を用いる。

1.2 研究の目的

本論文では、分散表現を用いて議論中での発言に含まれる単語の関連度を計算し、話題の変化を観測する手法を提案する。

1.3 本論文の構成

本論文の構成を以下に示す。2 章では要約手法に関する研究と、分散表現に関する先行研究を紹介する。次に、3 章では発言の要約手法の説明を行い、4 章では分散表現を用いた単語集合間の関連度計算について説明する。そして、5 章では話題

転換点の検出の評価実験について説明する．最後に 6 章で本論文のまとめと考察を示す．

第2章

関連研究

2.1 序言

本研究を構成する重要な要素として，COLLAGREE，議論支援，話題遷移検出，重み付け，分散表現 が挙げられる．本章では各要素について説明しながら関連研究を紹介する．

本章の構成を以下に示す．まず，2.2 節では，COLLAGREE の概要について簡単に述べる．2.3 節では，既存の議論支援に関する研究について述べる．2.4 節では話題に関する本研究の見方や既存の話題に関する研究について述べる．次に，2.5 節では本研究の重要な要素である重み付けについて述べる．2.6 節では本研究のもう 1 つの重要な要素である分散表現及び分散表現を使用した話題関連研究について述べる．最後に，2.7 節で本章のまとめを示す．

2.2 COLLAGREE

2.2.1 概要

COLLAGREE は掲示板のような議論プラットフォームをベースにしており，各ユーザーが自由なタイミングで意見を投稿，返信できる．こうした議論掲示板は基本的には 1 つの議論テーマに対して関連するテーマを扱った複数のスレッドから構成される．スレッドとはある特定の話題・論点に関する 1 つのまとまりを指す．例えば，図 2.1 のようにスレッドを立てたユーザーの発言が親意見となり，他のユーザーが子意見として親意見に返信し，子意見に対して孫意見が存在する場合もある．

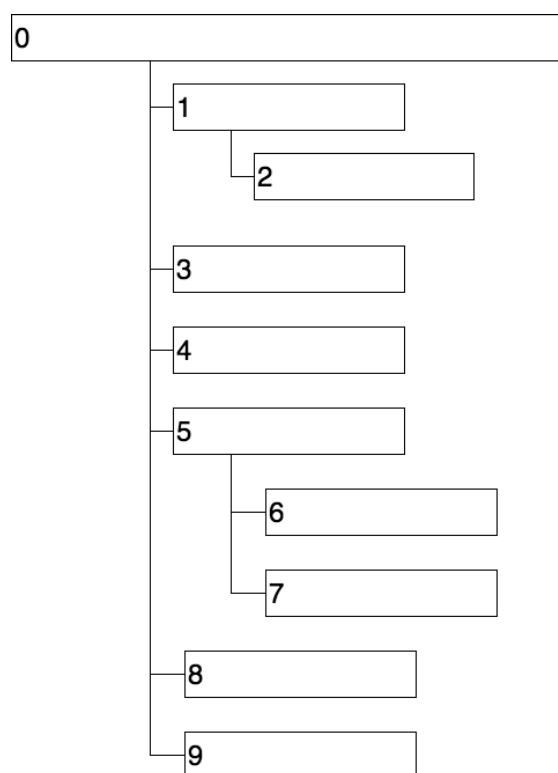


図 2.1: COLLAGREE のスレッドの例

しかし，スレッドや返信は必ずしも正しく使われるとは限らず，単なるチャット

と同様に扱われてしまいスレッドが乱立してしまうこともある。

2.2.2 ファシリテーター

COLLAGREE ではファシリテーターと呼ばれる人物が議論のマネジメントを行っている。ファシリテーターは”促進者”を意味し、議論そのものには参加せず、あくまで中立的な立場から活動の支援を行うようにし、自分の意見を述べたり自ら意思決定をすることはない。基本的な役割として”議論の内容の整理”，”議論の脱線防止”，”意見の促し”等が挙げられる。伊藤ら [15] によってファシリテーターの存在や手腕が合意形成に強く影響を与えることが示されている。

2.3 議論支援

議論は収束を目標として互いの意見を述べ合うことで、何かを決定する際において重要な社会活動である。故に注目が集まり、COLLAGREE 以外でも議論内容の把握支援を行い、議論の進行を支援するための研究が行われている。小谷ら [17] は好意的発言影響度を取り入れた議論支援システムを開発した。システムは議論中の発言の意図や内容に加えて、発言に対するリアクション (同意、非同意、意見) などから議論進行をモニタリングしている。モニタリングの結果を基にして議論の活性化や深化に対して参加者が果たしている役割を”好意的発言影響度”として定量化して表示する。

本研究はファシリテーターに対する支援を目的としており、一般参加者や学習者の議論活性化及び収束に向けた支援が目的ではない。

2.4 話題遷移検出

話題に関連する研究は長い間行われており，古くは「会話分析」という学問に始まる．会話分析は会話を始めとする相互行為の組織(構造)を明らかにしようとする社会学の研究分野であり，相互行為の分析法も取り扱う．1960年代に Harvey Sacks と Emanuel A. schegloff によって創案された．話題というものは、本来流れの一時点で簡単に区切ることのできるものではない．しかし、研究によっては話題内容と連鎖組織及び言語形式との関連についての分析を行うことがある．筒井[1]は話題を区切るのに次のように基準を立てた．

1. それまで話題となっていた対象や事態とは異なる，新しい対象や事態への言及
2. すでに言及された対象や事態の異なる側面への言及
3. すでに言及された対象や事態の異なる時間における様相への言及
4. すでに言及された対象や事態について，それと同種の対象や事態への言及
5. すでに言及された個別の対象や事態の一般化

本研究では上記の基準を参考に話題変化の判定の評価の基準を設けている．

話題の遷移に基づいた文章の分割は人間によるテキスト全体の内容の把握の容易化や複数のテキストに対する自動分類や検索の精度向上のために研究されている．以下では話題の遷移に関する関連研究について述べる．

2.4.1 テキストセグメンテーション

テキストセグメンテーションは複数のトピックが混合的に書かれている非構造的である文書をトピックに応じて分割する手法である。

TextTiling

Hearst ら [6] は複数の単語を連結した 2 つのブロックをテキストの初めから末尾まで動かしていきブロック間の類似度を計算する手法 (TextTiling) を提案した。類似度はブロック間で共通して使われる単語数などで計算する。図 2.2 に手法の簡単な流れを示す。

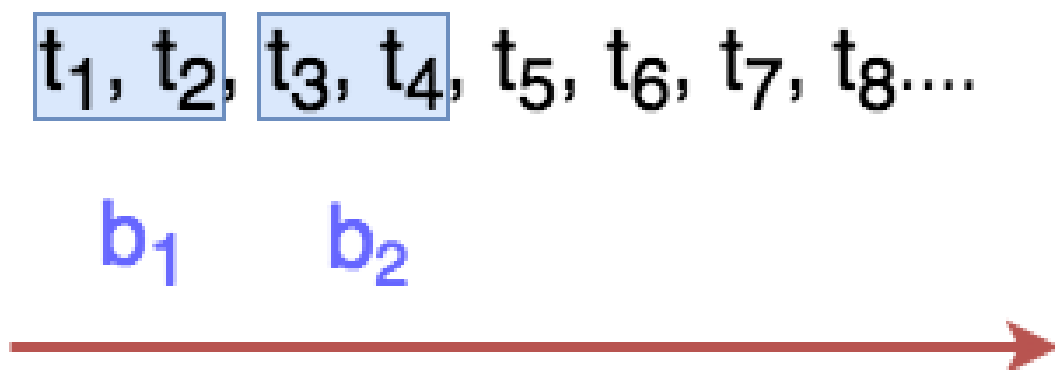


図 2.2: TextTiling の流れ

計算された類似度をグラフにすると図 2.3 のようになる。縦軸はブロック間の類似度、横軸はブロック間の番号を表し、グラフ中の番号は段落番号、垂直線は選出された文の境界位置を表す。グラフは各ブロック間の類似度を繋いだものと、更に平滑化したものの 2 つが描かれている。グラフの谷のような場所は左右のブロック間の類似度が下がる場所を表し、使われる単語が変わったこと、すなわちトピックが変わったことを意味し、谷の付近に境界を設けることで文書をトピックごとに区切ることを可能としている。

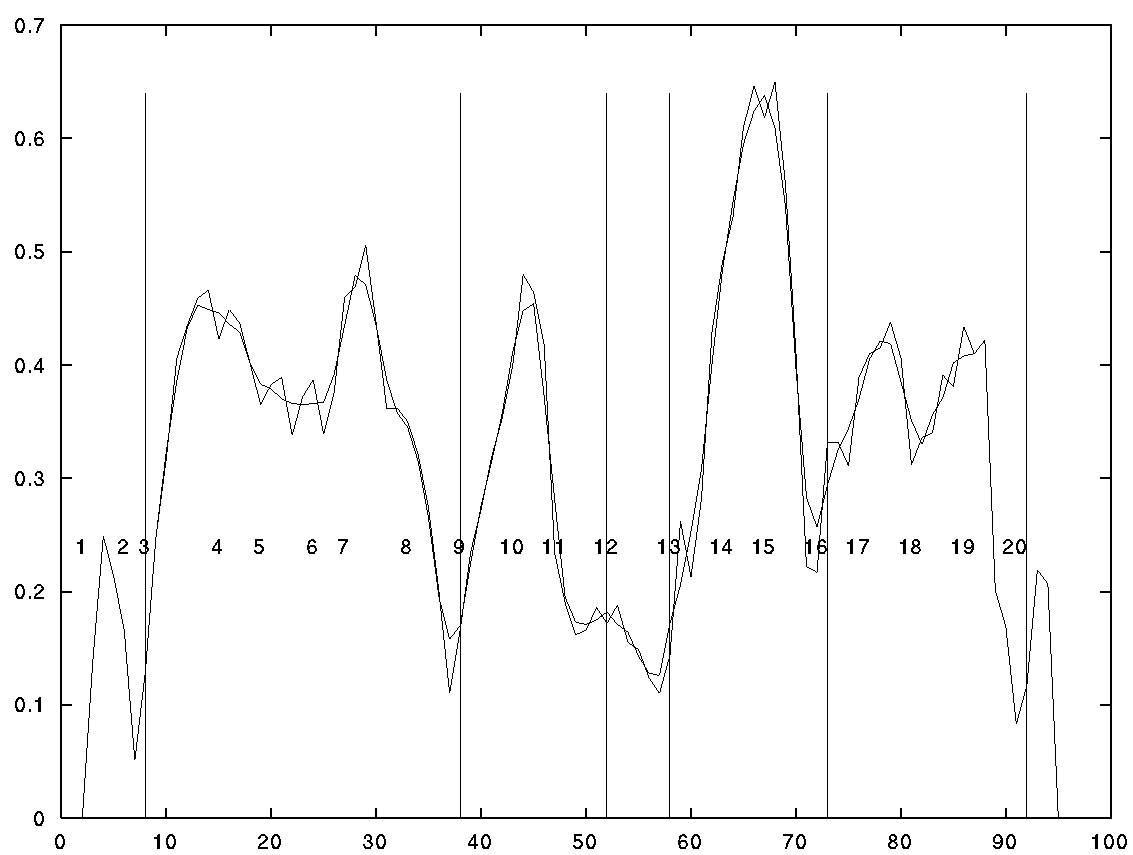


図 2.3: TextTiling のグラフ

テキストセグメンテーションを使用した関連研究

別所 [19] は単語の共起頻度行列を特異値分解で次元を削減して作成した単語の概念ベクトルを用いて新聞記事の分割を行っている。ブロック間の類似度を計算する際にブロック中の単語の内、自立語にのみ概念ベクトルを付与し左右のブロックの和ベクトル (または重心ベクトル) の余弦測度を求め、類似度 (または結束度) としている。

テキストセグメンテーションは基本的にある地点 A で話題に沿って分割をするか考える際に、地点 A より先の情報を使うことができる。すなわち、ある程度の文章が現れてから話題が変わったかどうかの判定をしており、リアルタイムでの動作を想定していない。本研究はリアルタイムでの動作を想定し、ある発言 B が話題の変化を起こすかどうか判定する際に B より先の情報を使うことなく判定している点で異なる。

2.4.2 トピックモデル

トピックモデルは、確率モデルの一種にあたり、文章中の「単語が出現する確率」を推定している。単語が出現する確率をうまく推定することができれば、似たような単語が出てくる文章が把握できる。すなわち、トピックモデルとは「文書における単語の出現確率」を推定するモデルといえる。

ユニグラムモデル

図は四角の箱が文書を表し、中身の色がトピックを表す。ユニグラムモデルは、すべてのボックスが同じ青色である。全ての文書の単語は、一つのトピックから生成されたものと仮定するモデルである。

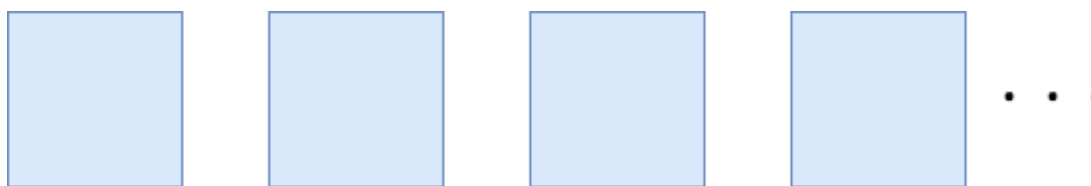


図 2.4: ユニグラムモデル

混合ユニグラムモデル

混合ユニグラムモデルは、各ボックスの色が異なっている。つまり、各文書に一つのトピックがあり、そのトピックから文書の単語が生成されると仮定するモデルである。



図 2.5: 混合ユニグラムモデル

LDA

LDA(Latent Dirichlet Allocation) では、ボックスの中で色が異なっている。つまり、各文書は複数のトピックで構成されていて、各トピックの単語分布を合算した形で単語が生成されていると仮定を行う。

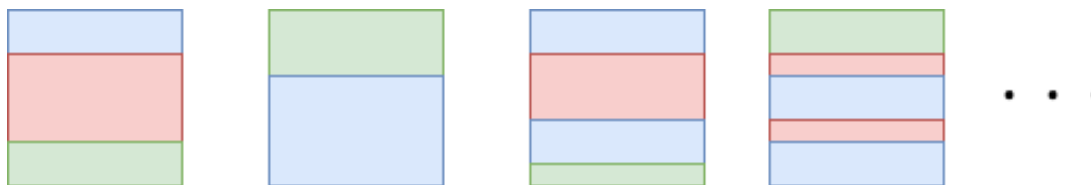


図 2.6: LDA

トピックモデルは基本的に学習の際にトピックの数を指定する必要がある、指

定に伴いトピック数が制限されてしまう．すなわち，トピックの数によっては変化に対応できない話題が存在することがあり得る．本研究は分散表現を用いており，トピック数の指定がない点で異なる．

2.5 重み付け

重み付けは情報検索を主とする分野で使われる方法で，蓄積された情報中の語の索引語，すなわち特定の情報の特徴を表し検索の手掛かりとなる語，としての重要度を数値的に表現し，それぞれの語の重要度に応じて重みを付け，集計して総合評価を出す手法である．出された総合評価はスコア等と呼ばれる．一般的に語の重要度は整数または実数値で与えられる．自動的な重み付けにおいては語の出現頻度情報や出現位置を利用して数値を与える．本研究で行う重み付けは下記の2種類の重み付けを基にしている．

2.5.1 Okapi BM25

Okapi BM25[5] は出現頻度情報を用いる重み付け手法の代表の1つである．Okapi BM25 ではTF(Term Frequency 単語の出現頻度)[12]，IDF(Inverse Document Frequency 逆文書頻度)[7]，DL(Document Length 文書長)の3つを用いて重要度の計算を行う．各用語について説明を行う．

TF

TF は単語の出現頻度を表し，文書中において出現頻度の高い単語は重要であるという考え方に基づく．ある単語 t_i の文書 D_j 中における出現頻度重み $tf_{i,j}$ は式 (2.1)

のようにして求められる.

$$tf_{i,j} = \frac{n_{ij}}{\sum_k n_{kj}} \quad (2.1)$$

ここで n_{ij} は文書 d_j における単語 t_i の出現回数, $\sum_k n_{kj}$ は文書 d_j におけるすべての単語の出現回数の和である.

IDF

IDF は逆文書頻度を表し, 多くの文書において出現頻度の高い単語は重要ではないという考え方に基づく. IDF は多くの文書に出現する語, すなわち一般的な語の重要度を下げ, 特定の文書にしか出現しない単語の重要度を上げる役割を果たす. ある単語 t_i の逆文書頻度重み idf_i は式 (2.2) のようにして求められる.

$$idf_i = \log \frac{|D|}{|\{d : d \ni t_i\}|} \quad (2.2)$$

ここで $|D|$ は総文書数, $|\{d : d \ni t_i\}|$ は単語 t_i を含む文書数である.

DL

DL は文書長を表し, ある単語の出現回数が同じ2つの文書について, 総単語数の少ない文書と多い文書では, 前者のほうがより価値があるという考え方に基づく. ある文書 d_j の文書長重み ndl_j は式 (2.3) のようにして求められる.

$$ndl_j = \frac{dl_j}{ave(dl)} \quad (2.3)$$

ここで dl_j は文書 d_j の総単語数, $ave(dl)$ はすべての文書の平均 dl を表す.

上記の3つの重みを用いて Okapi BM25 は (2.4) のように単語 t_i の文書 D_j 中にお

ける統合重み cw_{ij} を求める．

$$cw_{ij} = \frac{tf_{i,j} \cdot idf_i \cdot (k_1 + 1)}{k_1 \cdot (1 - b + b \cdot ndl_j) + tf_{i,j}} \quad (2.4)$$

ここで定数 k_1 と b について説明する． 2つの定数はどちらもチューニングの役割を果たすもので k_1 は単語の出現頻度による影響を， b は文書の長さによる影響を調節する．

2.5.2 LexRank

Erkan ら [3] によって考案された LexRank は Google の PageRank[14] を使用した文章要約アルゴリズムで， 文の類似度を計算して次の 2つの基準に基いて文の重要度を計算する．

1. 多くの文に類似する文は重要な文である．
2. 重要な文に類似する文は重要な文である．

LexRank でいう類似度は簡単にいえば 2文がどれだけ共通の単語を持つかということを表し， 文を TF-IDF を用いてベクトル化して Cosine を求めることで類似度としており， 式 2.5 に沿ってベクトル x と y の Cosine が計算される．

$$\text{idf-modified-cosine}(x, y) = \frac{\sum_{w \in x, y} tf_{w,x} tf_{w,y} (idf_w)^2}{\sqrt{\sum_{x_i \in x} (tf_{x_i,x} idf_{x_i})^2} \times \sqrt{\sum_{y_i \in y} (tf_{y_i,y} idf_{y_i})^2}} \quad (2.5)$$

文の間の Cosine 類似度をグラフとして可視化すると図 2.7 のようになる． 各エッジは文の間の Cosine 類似度を表し， $dXsY$ は文書 X の Y 番目の文を示す．

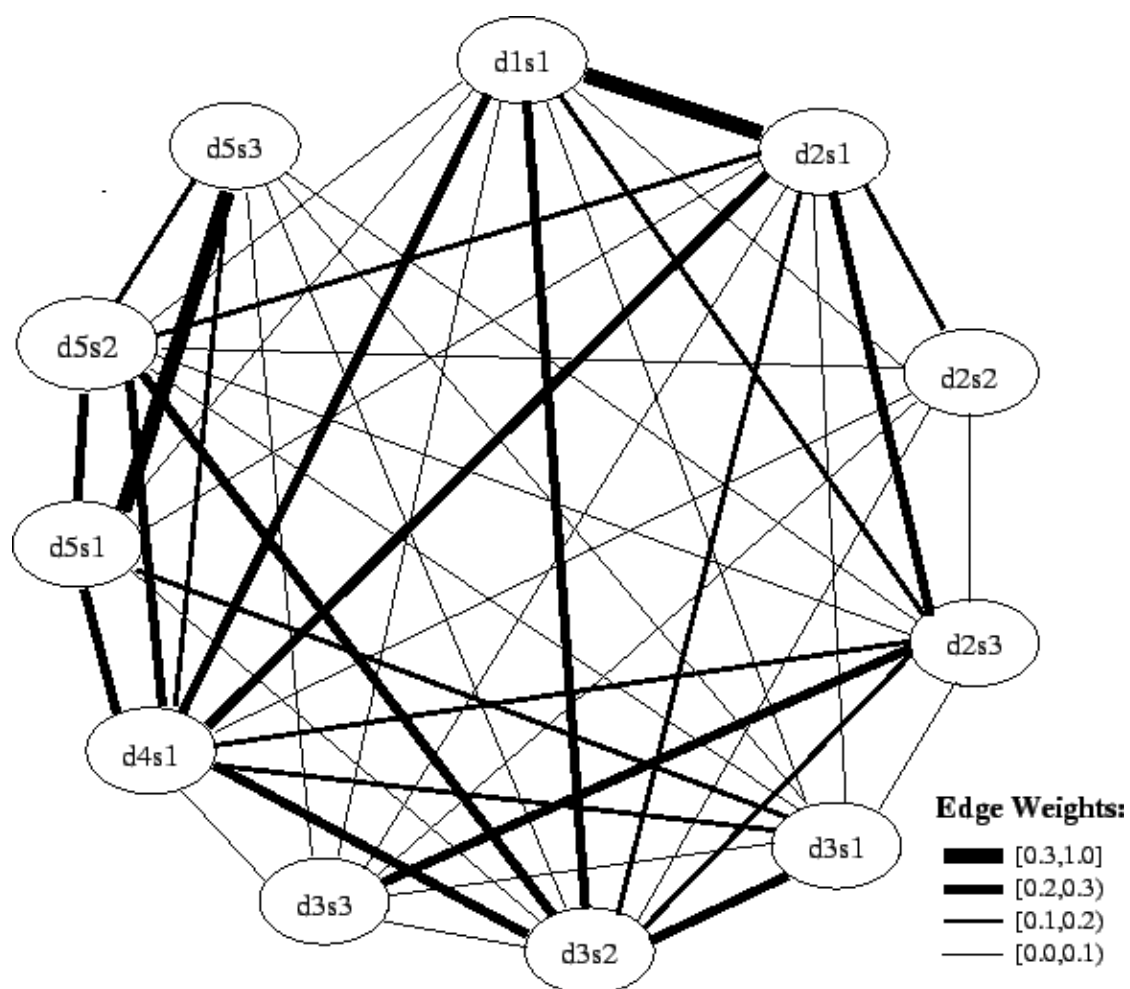


図 2.7: 類似度グラフの例

その後、Cosine 類似度が閾値を超えたかどうかを基に隣接行列が作成される。隣接行列はグラフを表現するために用いられる行列で、あるノード v と w の間にエッジの有無が行列の (v, w) 成分に割り当てられる。隣接行列の各要素を類似している文の数で割り、確率行列に変換した後、**Algorithm1** に従って行列の固有ベクトル \mathbf{p} が計算される。求められた固有ベクトルが LexRank スコアとなる。

Algorithm 1 ベキ乗法の計算アルゴリズム

```

1: Input : 確率的かつ既約かつ非周期的な行列  $M$ 
2: Input : 行列サイズ  $N$ , 誤差許容値  $\epsilon$ 
3: Output: 固有ベクトル  $\mathbf{p}$ 
4: procedure POWERMETHOD( $M, N, \epsilon$ )
5:    $\mathbf{p}_0 = \frac{1}{N} \mathbf{1}$ ;
6:    $t = 0$ ;
7:   repeat
8:      $t = t + 1$ ;
9:      $\mathbf{p}_t = M^T \mathbf{p}_{t-1}$ ;
10:     $\delta = \|\mathbf{p}_t - \mathbf{p}_{t-1}\|$ ;
11:  until  $\delta < \epsilon$ 
12:  return  $\mathbf{p}_t$ ;

```

LexRank スコアを計算する一連のアルゴリズムを **Algorithm2** に示す。7~14 行目では隣接行列が作成されており、15~18 行目では LexRank が計算されている。

2.6 分散表現

分散表現は単語を低次元または高次元の実数ベクトルで表現する技術で Harris[4] が提唱した”分布仮説”(”同じ文脈で出現する単語は類似した味を持つ傾向があり、単語はその単語とともに出現する単語等によって特徴づけられる。”という考え方)に基づく。

Algorithm 2 LexRank スコアの計算アルゴリズム

```

1: Input :  $n$  個の文からなる配列  $S$ , コサイン類似度の閾値  $threshold$ 
2: Output : 各文の LexRank スコアを格納した配列  $L$ 
3: Array  $CosineMatrix[n][n]$ ;
4: Array  $Degree[n]$ ;
5: Array  $L[n]$ ;
6: procedure LEXRANK( $S, t$ )
7:   for  $i \leftarrow 1$  to  $n$  do
8:     for  $j \leftarrow 1$  to  $n$  do
9:        $CosineMatrix[i][j] = \text{idf-modified-cosine}(S[i], S[j])$ ;
10:      if  $CosineMatrix[i][j] > threshold$  then
11:         $CosineMatrix[i][j] = 1$ ;
12:         $Degree[i] ++$ ;
13:      else
14:         $CosineMatrix[i][j] = 0$ ;
15:   for  $i \leftarrow 1$  to  $n$  do
16:     for  $j \leftarrow 1$  to  $n$  do
17:        $CosineMatrix[i][j] = CosineMatrix[i][j] / Degree[i]$ ;
18:    $L = \text{PowerMethod}(CosineMatrix, n, \epsilon)$ ;
19:   return  $L$ 

```

2.6.1 単語文脈行列

単語文脈行列は分散表現の最も基本的な形式で図 2.8 のような形の行列で表される。

単語の前後に出現する単語

	have	new	drink	bottle	ride	speed	read
beer	36	14	72	57	3	0	1
wine	108	14	92	86	0	1	2
car	578	284	3	2	37	44	3
train	291	94	3	0	72	43	2
book	841	201	0	0	2	1	338

コーパス
中の単語

図 2.8: 単語文脈行列

行列中の各要素 M_{ij} は単語 i と文脈 j の共起頻度を表しており、例として青い四角は train と ride が 72 回共起したことを表している。各行 M_i は単語 i の意味ベクトルを表し、例として赤い四角は”beer”の単語ベクトルを表している。

また、単語の類似度を式 2.6 のように単語の意味ベクトルのコサイン類似度で求めることができる。

$$\cos\theta = \frac{M_i \cdot M_j}{|M_i||M_j|} \quad (2.6)$$

式 2.6 を図 2.8 の行列に適応させると beer と wine のコサイン類似度は約 0.941, beer と train のコサイン類似度は約 0.387 となり, train よりも wine の方が beer に類似していることが分かる。

2.6.2 word2vec

2.6.1 節で説明した単語文脈行列から得られる単語ベクトルは単語の類似度を求めることは出来たが，他の数学的処理には対応していなかった．Mikolov ら [13] が開発した word2vec はニューラルネットワークを用いて分散表現の生成を行う手法で，CBOW と skip-gram の 2 種類の学習方法が存在する．図 2.9，図 refFig:skip-gram はそれぞれ CBOW と skip-gram の構造を表す．

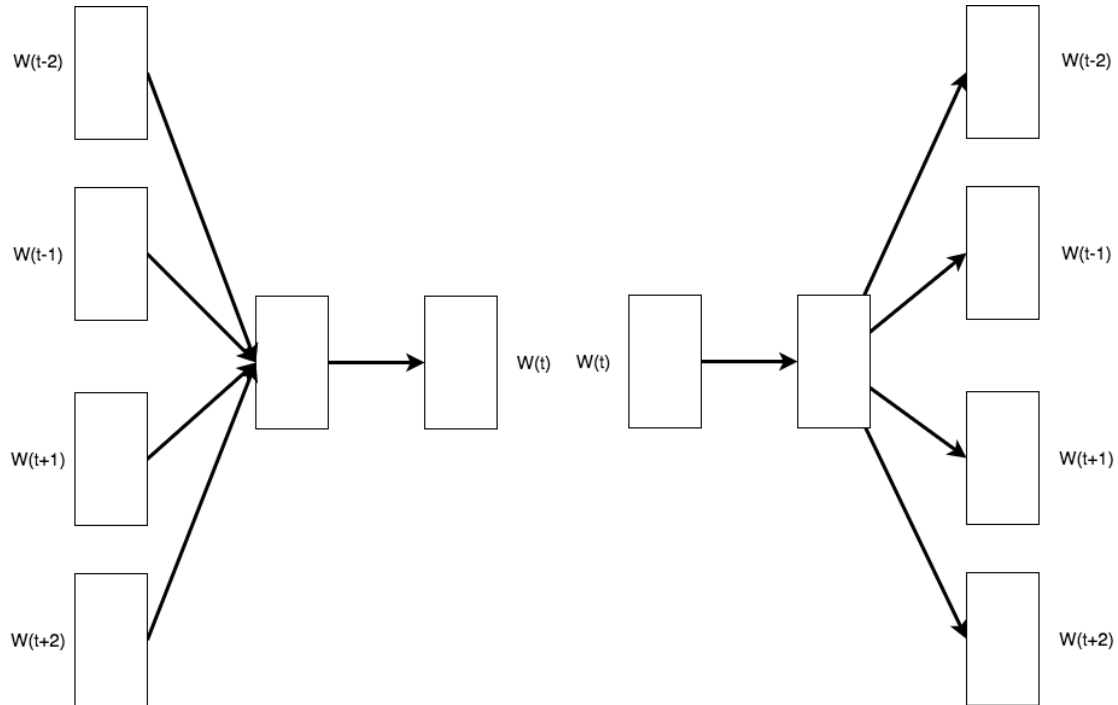


図 2.9: CBOW

図 2.10: skip-gram

CBOW は周辺の単語 $W(t-2) \cdots W(t+2)$ を入力として，現在の単語 $W(t)$ を予測することを目指して学習する．逆に skip-gram は現在の単語を入力として，周辺の単語を予測することを目指して学習する．どちらの手法でも中間層と単語の変換処理が行われており，学習によって得られる単語と中間層での数値が対応した行列が単語の分散表現モデルとなる．

word2vec が従来の手法と比べて大きく異なる点は ニューラルネットによる学習で単語の類似度の計算に加えて、ベクトルの加減算が単語の意味の加減算に対応しているということである。また、従来の手法を用いた単語ベクトルよりも類推精度が高いことが Levy ら [10] によって示されている。例えば、 $X = \text{vector}(\text{"biggest"}) - \text{vector}(\text{"big"}) + \text{vector}(\text{"small"})$ として計算されたベクトル X に最も類似したベクトルを探すことで biggest が big に類似しているのと同じ意味で small に類似している単語 smallest を見つけることができる。ただし、分散表現モデルが十分に訓練されていることが前提である。

2.6.3 fastText

fastText[2][8] は word2vec を発展させた手法でより大きな語彙や多くの稀な単語に対応することができ、学習の速度を上昇させることにも成功している。

fastText は skip-gram モデルを採用しており、学習の際に単語だけでなく部分語(単語を構成する文字のまとまり)についても考慮する。以下に単語 w_t が与えられた時に予測した文脈単語 w_c の間のスコア関数を示す。

$$s(w_t, w_c) = \mathbf{u}_{w_t}^T \mathbf{v}_{w_c} \quad (2.7)$$

$$s(w_t, w_c) = \sum_{g \in \mathcal{G}_w} \mathbf{z}_g^T \mathbf{v}_{w_c} \quad (2.8)$$

式 2.7 は fastText 以前の手法でのスコア関数を表し、式 2.8 は fastText でのスコア関数を表す。 \mathbf{u}_{w_t} , \mathbf{v}_{w_c} はそれぞれ単語 w_t , w_c を実数ベクトルで表したもので、式 2.8 において \mathcal{G}_w , \mathbf{z}_g はそれぞれ単語 w の n-gram の集合と n-gram g を実数ベク

トルで表現したものを表す．式 2.7 では単語と文脈単語の間のスカラー積をスコアとしているが，式 2.8 では単語の n-gram と文脈単語の間のスカラー積の合計をスコアとしている．式 2.8 の手法を用いることで従来のモデルでは考慮されていなかった”活用形”を考慮できるようになった．例として，単語 go と goes と going は全て go の活用形であるが字面は異なるので従来のモデルでは異なる単語として学習されていたが，fastText では部分語である”go”を 3 つ全てで学習することで意味の近い単語として学習することが可能となることが挙げられる．

本研究では分散表現の手法として fastText を使用している．

2.6.4 分散表現を使用した話題関連研究

分散表現は関連語を導出できることから分散表現を話題関連に用いる研究が行われている．中野ら [18] は分散表現を用いて雑談対話システムでのシステム側の応答生成を行っている．図 2.11 に話題展開システムの構造を示す．

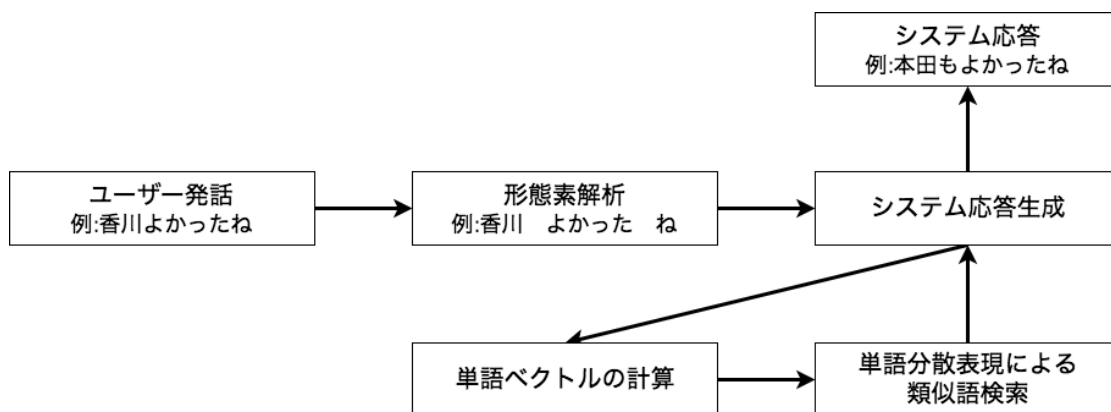


図 2.11: 話題展開システムの構造

システムではユーザ発話を形態素解析して検出された単語を単語分散表現による類似語検索から得られた結果と入れ替えることでシステム応答の生成を行って

いる。

また, Li ら [11] は分散表現を用いて Twitter のツイートにトピックカテゴリに分類する分類器 TweetSift を提案している. 図 2.11 にトピックカテゴリの予測のワークフローを示す.

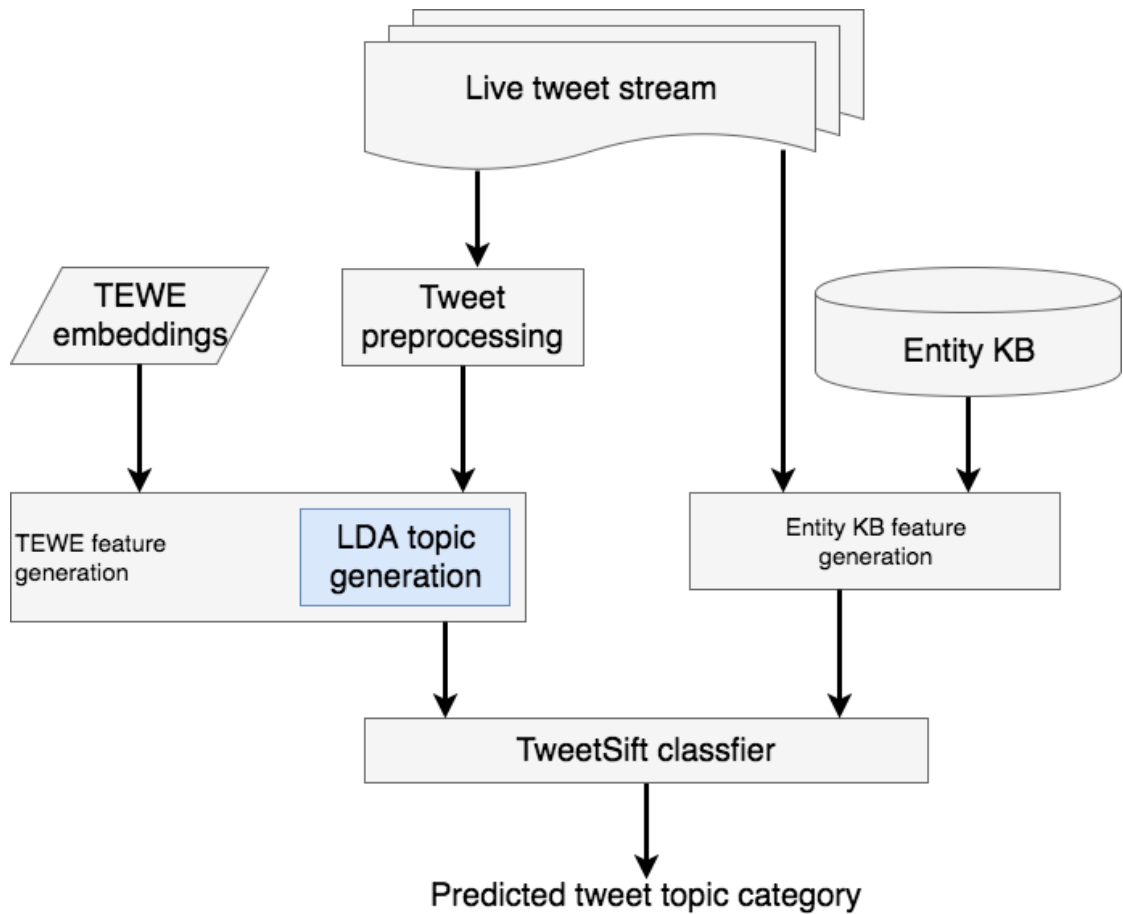


図 2.12: TweetSift によるトピック予測のワークフロー

図の右のフローではツイートデータとスクレイピングで作成された知識ベースを用いて知識ベース特徴を生成している. 図の左のフローでは前処理を行ったツイートと作成済み分散表現モデルを用いて分散表現特徴を生成している. 2つの特徴を用いて SMO(Sequential Minimal Optimization) によってトピックが予測されている. Li らの研究で用いられる分散表現は学習の際に単語だけでなく, LDA

を用いて予測した単語のトピックも使用されている。

本研究は Web 上での議論を対象としており，対話や SNS を対象としていない点で異なる。また，中野らの研究とは文の生成を行わない点でも異なり，Li らの研究とは分散表現の学習の際にトピックを用いていない点でも異なる。

2.7 結言

本章では，COLLAGREE や他の議論プラットフォームでの議論支援研究を紹介し，本研究との違いを説明した，また，本研究と類似した研究や本研究での重要な要素，及び要素を使用した関連研究についても説明し，本研究の立ち位置を明らかにした。

第3章

通知システム

3.1 序言

本章では通知システムの概要について説明する．以下に本章の構成を示す．まず3.2節でシステム全体の動作の流れを示し，アルゴリズムについても説明を行う．次に3.3節ではシステムの詳細を説明する．システムで扱う発言データの形式や発言間の類似度計算について述べる．最後に，3.4節で本章のまとめを示す．

3.2 システムの動作の流れ

擬似コードを **Algorithm3** に示し，図示したものを図3.1に示すシステムの動作の流れについて説明する．発言 R が投稿された時，過去に投稿された発言と類似度の計算を行い類似度が閾値を超えていれば2つの発言が同じ話題であるとみなし，発言 R と同じ話題である発言の集合 SG に登録する．作業を繰り返し全ての発言との計算が終了した後， SG が空集合である，すなわち発言 R と同じ話題である発言がない場合に話題を変化させる発言であると判定して通知を行う．

Algorithm 3 システムの流れ

```

1: Input : 発言  $R$ 
2: Output : 通知判定  $Notify$ 
3:  $PG$  = 過去の発言の集合;
4: procedure TOPICCHANGE( $R$ )
5:    $SG = \{\}$ ;
6:   for Each  $pastR \in PG$  do
7:      $sim = \text{similarity}(R, pastR)$ 
8:     if  $sim > \text{threshold}$  then
9:        $SG.append(pastR)$ 
10:   $Notify = \text{False}$ 
11:  if  $SG == \{\}$  then
12:     $Notify = \text{True}$ 
13:  return  $Notify$ 

```

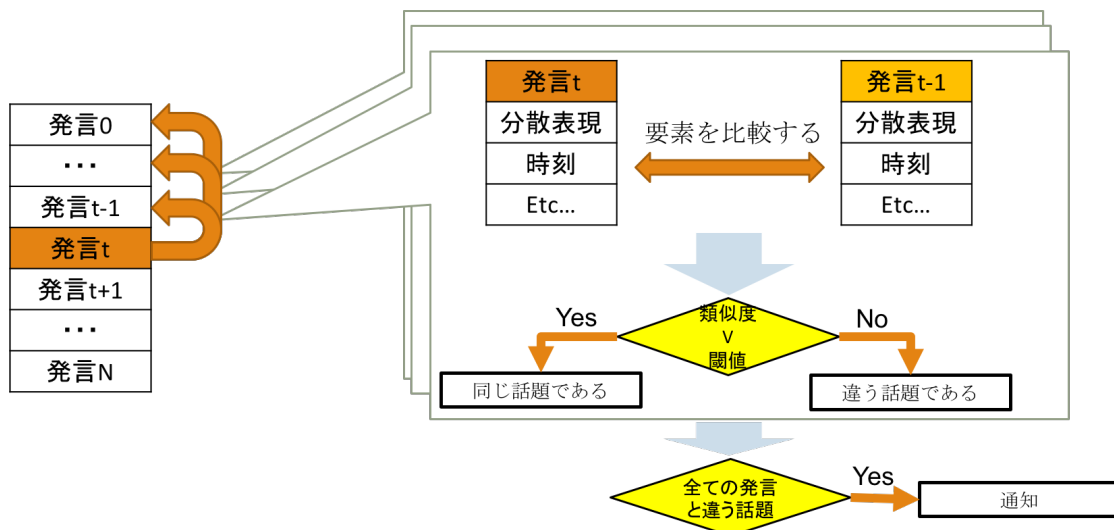


図 3.1: システムの流れ

3.3 システム詳細

3.3.1 発言データ

本システムで扱う”発言データ”は単なる文字列ではなく、タイムスタンプ等の他のデータを持つリスト形式のデータである。表 3.1 にデータの一部を示す。また、以下で本研究で使用する発言データの要素について説明する。

① id

発言データを識別するための番号で、全て整数値で表される。

② title

スレッドのタイトル名を表す文字列で、発言がスレッドの先頭でない限りは NULL となる。

③ body

発言の内容を表す文字列。

④ parent-id

発言の返信先、すなわち親発言の id を表す番号で、親発言がある場合は親発言の id と同じ番号になり、ない場合は NULL となる。

⑤ created-at

発言が投稿された時間を示すタイムスタンプ。

id	title	body	parent-id	user-id	(以降は省略する)
18	オンラインでの議論に関する実験	参加者の皆様が集まるまでお待ち下さい。	NULL	1	(以降は省略する)

表 3.1: 発言データ

3.3.2 発言間の類似度計算

発言間の類似度計算は次の2段階で行われる。

① 文章間の類似度計算

発言データ中の title と body, すなわち発言の内容の類似度の計算を行う。title が NULL でない場合は title と body を改行コードで繋いで1つの文章とする。

文章間の類似度計算の手法については4章で詳しく述べる。

② 総合類似度計算

上記の①で計算された発言の文章間の類似度に発言間の時間差と返信関係を組み合わせることで総合類似度を求める。

時間差評価値

発言 *new* と以前の発言 *old* 間の時間差を式 3.1 に基いて正規化された評価値として求める。

$$tValue = 1 - \frac{epoch(new.created) - epoch(old.created)}{maxTime} \quad (3.1)$$

ここで関数 *epoch*, 定数 *maxTime* について説明する。*epoch* は与えられたタイムスタンプをエポック秒に変換する。*maxTime* は議論の制限時間を表す。時間差評

価値は 2 発言間の時間差が小さいほど関連が強いとみなし, 0 から 1 に近づく.

返信距離

発言 *new* と以前の発言 *old* 間の返信距離を Algorithm4 に基いて再帰的に求める.

Algorithm 4 返信距離

```

1: Input : 発言 new, 発言 old, 返信距離 dist                                ▷ 初期値 dist=1
2: Output : 返信距離 dist
3: PG = ID に対応づけられた過去の発言の集合;
4: procedure REPLYDIST(new, old, dist)
5:   if new.parent-id == NULL then
6:     return 0
7:   else if new.parent-id==old.id then
8:     return dist
9:   else
10:    parent = PG[new.parent-id]
11:    dist + = 1
12:    return REPLYDIST(parent, old, dist)

```

7 ~ 8 行目, 9 ~ 12 行目で示すように発言の id が一致した場合は現在の返信距離を返し, 一致しなかった場合は返信距離を 1 増やして *new* の親発言 *parent* と *old* の返信関係を再帰呼び出しで求め, 戻り値を返す. また, 5 ~ 6 行目で示すように 2 発言間が返信関係になかった場合は 0 を返す.

総合類似度

総合類似度は前述の発言内容の類似度, 時間差評価値, 返信距離によって求められる. 返信距離が 0 である時, すなわち 2 発言が異なるスレッドに属している場合は類似度と時間差評価値から総合類似度を計算する. 類似度だけでなく時間差評価値を使用するのは, 総合類似度だけで判断してしまうと議論の終盤になって発言数が多くなってきた時に新しく投稿された発言が多くの古い発言と類似している

と判断されてしまうことがあり得るからである．議論は基本的に少し前の発言に関連して行われることが多いことから時間差評価値を使用して時間的に近いもののほど総合類似度が上昇するようにする．具体的には式 3.2 のように計算される．

$$tSim = tValue * tWeight + sim * (1 - tWeight) \quad (3.2)$$

ここで変数 sim ，定数 $tWeight$ について説明する． sim は発言内容の類似度を表し，0 から 1 の値を取る． $tWeight$ は時間差評価値の総合類似度の計算における重要性を表し，0 から 1 の値を取る．また，返信距離が 0 でない，すなわち 2 発言が同じスレッドに属している場合は何らかの関連があると考えられることから，発言内容の類似度を総合類似度とする．

3.4 結言

本章では通知システムの動作の流れやアルゴリズムについて説明し，扱うデータの形式や内容，及び発言内容の類似度を除く発言間の類似度計算の手法について説明した．また，時間的に近い議論ほど類似度が上昇するように時間差評価値と発言内容の類似度の 2 つを用いて総合類似度計算することを示した．

第4章

発言内容の類似度計算

4.1 序言

本章では発言内容，すなわち文字列の意味的類似度を計算する手法を提案する．以下に本章の構成を示す．4.2 節では類似度の精度を上昇させるために行う前処理について説明する．4.3 節では発言内容の類似度計算の手法について述べる．4.4 節では本章のまとめを示す．

4.2 前処理

分散表現による類似度計算で精度を上昇させるためには発言内容から余分な単語を取り除きいて重要な単語を抽出する，または極めて短く要約することが重要である．提案手法では形態素解析エンジン MeCab と 2.5 節で説明した okapiBM25 と LexRank を用いて発言の文章から重要単語を抽出し，抽出した単語の類似度を分散表現を用いて計算する．

4.2.1 MeCab

MeCab(めかぶ)[9] は京都大学情報学研究科-日本電信電話株式会社コミュニケーション科学基礎研究所共同研究ユニットプロジェクトを通じて開発されたオープンソース形態素解析エンジンである。MeCab に対して「MeCab はオープンソース形態素解析エンジンである。」と入力した際の結果を図 4.1 に示す。

```
Mecabはオープンソース形態素解析エンジンである。
Mecab 名詞,固有名詞,組織,*,*,*,*
は 助詞,係助詞,*,*,*,*,は,ハ,ワ
オープン 名詞,サ変接続,*,*,*,*,オープン,オープン,オープン
ソース 名詞,一般,*,*,*,*,ソース,ソース,ソース
形態素 名詞,一般,*,*,*,*,形態素,ケイタイソ,ケイタイソ
解析 名詞,サ変接続,*,*,*,*,解析,カイセキ,カイセキ
エンジン 名詞,一般,*,*,*,*,エンジン,エンジン,エンジン
で 助動詞,*,*,*,特殊・ダ,連用形,だ,デ,デ
ある 助動詞,*,*,*,五段・ラ行アル,基本形,ある,アル,アル
。 記号,句点,*,*,*,*,。,。,。
EOS
```

図 4.1: 解析結果

出力フォーマットは次の形式となっている。

表層形 \t 品詞, 品詞細分類 1, 品詞細分類 2, 品詞細分類 3, 活用型, 活用形, 原形, 読み, 発音

「読み」と「発音」は図 4.1 の”MeCab”のように不明であるものには付与されない。MeCab による形態素解析の結果、次の条件を満たす単語を除外している。

1. 品詞細分類に「数」を含む
2. 「読み」,「発音」が不明である
3. 品詞が「助詞」,「助動詞」,「記号」,「連体詞」のどれかである。
4. 1 文字のひらがなである

5. 品詞細分類に「接尾」または「非自立」を含む

上記の条件を満たす単語を除外したのは4.2.2節で説明する重み付けにおいて重要な単語であると判定されやすいが，??節で説明する類似度計算において精度を下げてしまうからである．単語の除外は重み付けにおいて文章を単語に分割する際に行われる．

4.2.2 重み付け

提案手法では2.5節で説明した okapiBM25 と LexRank の2種類の重み付け手法をを統合して発言の内容の文字列 *remark* 中の単語に対して重み付けを行う．アルゴリズムを **Algorithm5** に示す．

Algorithm 5 統合重みの計算アルゴリズム

```

1: Input : remark 発言内容の文字列
2: Output : combinedWeight remark 中の単語と重みを対応付けた連想配列
3: Array sentList; ▷ 以前に重み付けを行った最大 n 個前までの文章のリスト
4: procedure CALCCOMBINEDWEIGHT(reamrk)
5:   bm25Weight = calcBM25Weight(remark) ▷ 単語と重みの連想配列
6:   for Each sent ∈ remark do ▷ remark を句点, 改行コードで分割する
7:     sentList.append(sent)
8:   lexWeight = calcLexRank(sentList)
9:   for Each word ∈ bm25Weight.keys() do
10:     wordWeight = bm25Weight[word]
11:     if word is 固有名詞 then
12:       wordWeight *=2
13:     sentWeight = 0
14:     for Each sent ∈ remark do
15:       if word in sent then
16:         sentWeight += lexWeight[sent]
17:     combinedWeight[word] = wordWeight * sentWeight
18:   return combinedWeight

```

固有名詞は文章の中で重要な役割を果たす可能性が大きいと考え、12行目では固有名詞の単語重みを倍にしている。そして、14～17行目では word を含む全文章の重みの合計を求め、okapiBM25 による単語重みを掛け合わせたものを word の統合重みとしている。単語重みに単語を含む文章の重みを掛け合わせることで感嘆文のような文章そのものは重要でないが頻度の少ない単語を使用する文章中の単語が選ばれる可能性を下げている。

4.3 類似度計算

4.3.1 単語抽出

4.2 節で計算された単語重みの値が大きいものの上位 n 個までの単語を発言文章 remark において重要度の高い単語であるとして抽出する。単語重みが等しいものが複数あった場合は単語を昇順に並び替えて順序を付けている。また、使用する分散表現モデルに登録されていない単語は除外している。

4.3.2 分散表現による類似度計算

4.3.1 節で述べた手法を用いて 2 発言それぞれから抽出した単語集合の類似度を分散表現を用いて求める。それぞれの単語集合の単語ベクトルの平均を求め、2.6 の図 2.6 で述べたように Cosine 類似度を 2 平均ベクトル間で取っている。

4.4 結言

本章では発言内容の類似度を計算する手法について説明した。文章を単語に分割する手法と使用したツールと単語を除外する前処理についても述べた。また，okapiBM25 と LexRank を組み合わせた発言中の単語の重み付け手法，及び重み付けによって抽出した単語集合の類似度計算手法についても説明した。

第5章

評価実験

5.1 序言

本章では，COLLAGREE で行われた議論のデータを対象にした提案手法の評価実験について述べる．評価実験では同じテーマのもと行われた複数の議論を用意し，提案手法の有用性を示す．以下に本章の構成を述べる．

5.2 対象データ

5.2.1 議論データ

議論データは COLLAGREE 上で行われた別の実験での議論のものを使用する。
データの概要を以下に示す。

【実験概要】

グループ人数 : 2~3 名

議論時間 : 90 分前後

議論テーマ : 外国人観光客向けの日本旅行プランの決定

議論テーマ説明文 : みなさまに、外国人観光客向けの日本旅行プランを立てて
いただきます。 想定される旅行者の条件は以下の通りです。

- 英語は話せるが、日本語は話せない
- 初めての日本旅行である
- 日程は 6 泊 7 日
- ホテルは自分たちで手配できる
- 旅行のために貯金したので、金銭的には余裕があり、国内をいろいろとまわることが可能である
- 来日、帰国の際の空港は、どこでもかまわない
- 2 つのプランを比較したいと考えている（プランは 2 つ用意してください）

ファシリテータ : あり

5.2.2 評価データ

5.2.1 節で説明した議論データに対し，次に述べる基準でアノテーションを行ってもらった．基準を満たすと思われる発言に”1”のタグを，満たすと思われない発言に”0”のタグを付ける．

① それまで話題となっていた対象や事態とは異なる，新しい対象や事態への言及する発言

話されている内容が，以前と全く異なる対象や事態へと移行する位置でデータを区切る．

例 1:

(今までの話題:パック旅行はなぜ安いのかについて)

- A:ホテルが宿泊費の一部を出しているから安いのかな？
- B:おそらく。
- A:なるほど。
- B:沖縄行きも安いね。(今まで沖縄の話はされておらず，この後“沖縄行きのパック旅行”に話題が変わる(かもしれない))

例 2:

(今までの話題:外国人のツアー旅行の行き先について)

- A:他は寄らなくてもよいですか？(新しい行き先が出るように仕向けている)

② 既に言及された対象や事態の異なる側面への言及する発言

既に話題として取り上げられることについて、以前とは異なる側面から言及がなされる位置で区切る。

例 3:

(今までの話題:外国人のツアー旅行の行き先について)

- A:広島、長崎はどう？
- B:外国人観光客とか広島、長崎で見かけた覚えがない。
- A:ツアーに英語を話せるスタッフとか付けられるかな？(“ツアー旅行のスタッフ”に話題が変わる(かもしれない))

③ 議論のフェーズを移行させる(かもしれない) 発言

議論のフェーズを今までから移行させる(と思われる) 発言の位置で区切る。

例 4:

(今までの話題:外国人のツアー旅行の行き先について)

- A:八坂神社や清水寺など有名どころがたくさんありますし、魅力的だと思います
- B:京都周辺ツアー清水寺、金閣寺、銀閣寺、伏見稲荷大社、嵐山、など日本の建物や食べ物など広島長崎ツアー広島、長崎の戦争の地を見る事と、それぞれの場所で食べ物建造物を見るツアー(地名を挙げる段階から、各地点を結ぶツアープランへの作成段階に話題が変わる(かもしれない).)

例 5:

(今までの話題:外国人のツアー旅行のプランについて)

- A:京都周辺ツアー清水寺、清水焼体験、抹茶・和菓子など体験、きもの体験、金閣寺、嵐山、伏見稲荷大社その中で乗れそうなら屋形船などはどうでしょうか?
- B:屋形船、風情があって良いと思います。
- (途中省略)
- C: まとめると、・京都周辺ツアー京都周辺（八坂神社、清水寺、金閣寺、銀閣寺、伏見稲荷大社、嵐山、有馬温泉）、おいしい料理（豆腐など）、温泉、6泊7日ツアー
・広島長崎ツアー広島（3日）：広島原爆ドーム、平和記念公園、厳島神社、もみじまんじゅう、牡蠣、広島筆（メイクや書道なので使用する）、
お好み焼き、呉の戦艦、アナゴ（移動1日）長崎（3日）：ハウステンボス、
グラバー園、眼鏡橋、大浦天主堂、軍艦島、長崎ちゃんぽん、佐世保バーガー
この2プランで問題ないでしょうか？(初めて、2つのツアーの内容をまとめ、議論の収束に近づけた。)

また、ファシリテーターによる議論をコントロールするような発言も含む。

例 6:

- F:もし現在の旅先候補でよろしければ、具体的なプランづくりに移行したい
と思います。よろしいでしょうか？
- F:残り 20 分を切りました。皆様、いかがでしょうか？

以上の基準に沿ってタグを付けてもらい，“1”のタグが過半数以上付けられた発言を正解値=1，他を正解値=0 とした．

5.3 実験設定

5.3.1 パラメーター

本実験ではパラメーターは次の通りに設定した．前処理にて用いる okapiBM25 のパラメーターは $k1=2$, $b=0.75$ とし，LexRank のパラメーターは $n=50$, $threshold=0.7$ とした．分散表現として用いる fastText は次元数を 100 次元とし，学習データには wikipedia ダンプデータを用いた．総合類似度の計算に用いるパラメーターは $maxTime=5400(90 \text{ 分})$, $tWeight = 0.5$ とし，総合類似度の閾値は 0.8 とした．表 5.1 に実験の設定をまとめる．

5.3.2 比較手法

① 常時通知

最も単純かつ分かりやすい比較手法として，発言の内容に関係なく常に通知を行う手法を用いる．

okapiBM25	k1	2
	b	0.75
LexRank	n	50
	threshold	0.7
fastText	次元	10
	学習データ	wikipedia ダンプデータ
maxTime		5400
tWeight		0.5
類似度閾値		0.8

表 5.1: パラメーターの設定

② TF-IDF ベクトル

単語の意味は考慮せず出現頻度に基づく比較手法として、分散表現の代わりに TF-IDF で発言をベクトル化する手法を用いる。Algorithm5 の 5 行目で okapiBM25 の代わりに TF-IDF を用いて連想配列を求め、重みのベクトルに変換する。発言内容の類似度計算は提案手法と同じで Cosine 類似度を用い、以降も同じである。

5.3.3 評価指標

本実験では評価指標として適合率 (Precision), 再現率 (Recall), F 値 (F-measure) の 3 種類の指標を用いる。

適合率, 再現率, F 値はそれぞれ次のようにして求める。まず, 発言の通知を行うと判定した時を予測値=1, 通知を行わないと判定した時を予測値=0 とおく。次に, 予測値=1 かつ正解値=1 であるものの個数を *hits*(的中数), 予測値=0 かつ正解値=1 であるものの個数を *misses*(見逃し数), 予測値=1 かつ正解値=0 であるものの個数を *falseAlarms*(誤警報数) として数える。そして, 式 5.1, 式 5.2 及び式

5.3 に従って適合率，再現率，F 値を計算する．

$$Precision = \frac{hits}{hits + falseAlarms} \quad (5.1)$$

$$Recall = \frac{hits}{hits + misses} \quad (5.2)$$

$$F - measure = \frac{2 \cdot Precision \cdot Recall}{Precision + Recall} \quad (5.3)$$

3 つの値はどれも値が高いほど判定精度が高いことを示す．

5.4 実験結果

実験結果を表 5.2 に示す．

手法	平均評価指標		
	Pre	Rec	F
比較手法 1	0.3	0.5	0.4
比較手法 2	0.3	0.5	0.4
比較手法 3	0.3	0.5	0.4

表 5.2: 実験結果

5.5 結言

本章では本研究で提案する話題変化の判定手法が有用であることを実験により確認した．COLLAGREE にて行われた議論データに対して基準を満たすと思われるものにタグを付けてもらい評価データとした．評価実験では発言の内容に関係

なく常に通知する手法と分散表現の代わりに TF-IDF を用いて発言をベクトル化する手法を比較手法として用いた.

実験の結果, 提案手法は適合率と再現率の両方でバランス良く高い結果を示すことが分かり, 比較手法よりも良い結果を出すことがわかった.

第6章

結論

6.1 序言

6.2 今後の課題・展望

6.3 本研究のまとめ

謝辞

本論文を作成するにあたって、数多くの方々の御支援、御協力を頂きました。ここに、その方々に心から感謝の気持ちを申し上げます。

指導教官である名古屋工業大学大学院産業戦略工学専攻 伊藤孝行教授 に感謝致します。伊藤教授には、日々の研究活動やプレゼンテーションに対する御指導だけでなく、普段の生活に至るまで幅広い御指導を頂き、企業との共同研究という貴重な体験をさせて頂きました。また、学部生にもかかわらず3か月間という長い間の研究留学をさせて頂き、海外の研究生の研究に対する姿勢や考え方に関して多くを学べ、自身の研究に対する価値観を変える大変大きな経験をさせて頂きました。さらに、伊藤教授の熱心な御指導により、国際学会を含めた様々な場で発表させて頂く機会を得ることができ、自身の成長につながりました。心より感謝申し上げます。

留学中の指導教官であるシドニー工科大学 Quantum Computation & Intelligent Systems の Ivor Tsang 准教授 に感謝致します。Ivor 准教授には研究活動だけでなく、将来の進路や進学に至るまで幅広い内容に対して熱心な御指導を頂きました。Ivor 准教授から頂いた様々な助言が、私の研究成果につながりました。心より感謝申し上げます。

名古屋工業大学 Rafik Hadfi 特任助教 に感謝致します。Rafik 特任助教には、研

究活動やプレゼンテーションに対する御指導を頂きました。また、英語での発表や論文執筆に対する御指導を頂き、自身の成長につながりました。さらに、将来の進路に関する相談にも親身に乗って頂き、多くの心強い助言を頂きました。ここに感謝の意を表します。

香港科技大学大学院 Department of Computer Science and Engineering の Xingjian Shi 氏に感謝致します。Xingjian 氏には、この論文を書くことになったきっかけを頂けただけでなく、見ず知らずの学部生相手に 40 通近くにも及ぶメールでのやり取りにおいて既存手法に関する詳細な説明を頂き、ソースコードを提供いただいた上に、提案手法に関する助言も頂きました。ここに感謝の意を表します。

シドニー工科大学の Han Bo 氏, Donna Xu 氏には、留学期間中に何度も機械学習に関する相談に乗って頂き、研究に関する助言を頂いただけでなく、論文執筆、学会・ジャーナルへの投稿、さらには将来の進路に至るまで幅広い助言を頂きました。ここに感謝の意を表します。

共同研究先の株式会社ウェザーサービス様には、評価実験のデータ収集に御協力頂いただけでなく、気象学の視点から本研究で提案したモデルや評価実験について参考になる助言を頂きました。ここに感謝の意を表します。

共同研究先の株式会社 NEC ソリューションイノベータ 加藤憲昭氏 には、研究活動に対する御支援をして頂いた上に、本研究の今後の展望について参考になるご意見を頂きました。ここに感謝の意を表します。

名古屋工業大学伊藤孝行研究室の秘書である 杉山順子氏 には、研究室での事務業務など、学生たちがよりよい環境で研究を行えるための御支援を頂きました。

ここに感謝の意を表します。

名古屋工業大学工学部情報工学専攻 伊藤孝行研究室の先輩である 徳田渉先輩には、研究活動に対する御指導を頂き、研究に関して幾度となく相談に乗って頂きました。また、御自身の研究が忙しいにもかかわらず、留学中に論文執筆の指導や研究の協力をしていただき、大変支えられました。ここに感謝の意を表します。

伊藤孝行研究室の卒業生である 佐藤元紀先輩には、深層学習に関してわからないことがあった際に、幾度となく質問に答えていただきました。また、卒業されているにもかかわらず、深層学習に関する輪講会や勉強会に同伴して頂き、御指導を頂きました。ここに感謝の意を表します。

名古屋工業大学工学部情報工学専攻 伊藤孝行研究室の先輩である 早川浩平先輩には、研究活動に対する御指導を頂いた上に、研究室でより有意義に研究ができるよう取り計い頂きました。早川先輩の御指導、御力添えなくして順当な研究活動は行えなかったと思います。ここに感謝の意を表します。

名古屋工業大学工学部情報工学専攻 伊藤孝行研究室の先輩である 森顕之先輩には、研究活動やプレゼンテーションに対する熱心な御指導を頂き、自身の成長につながりました。また、先輩の研究に対する姿勢や言動から、論理的に話すことの大切さを学びました。ここに感謝の意を表します。

名古屋工業大学工学部情報工学科 伊藤孝行研究室の高橋一将君、仙石晃久君、Gu Wen 君、石田健太君、稲本琢磨君には同じ研究室の仲間として何度も助けられました。そしてこの伊藤研究室での有意義な時間を共に過ごすことができました。ここに感謝の意を表します。

また，友人の皆さんには貴重な時間と数々のご意見を頂きました．皆さんと過ごした時間はこれからも自身の励みとなると思います．

最後に自分の日々の生活を支えて頂いた家族に心より深く感謝いたします．

伊藤孝行研究室にて

2016 年 春

林政行

参考文献

- [1]
- [2] et al Bojanowski, Piotr. "enriching word vectors with subword information". *arXiv preprint arXiv:1607.04606*, 2016.
- [3] Gunes Erkan and Dragomir R. Radev. "lexrank: Graph-based lexical centrality as salience in text summarization". *Journal of Artificial Intelligence Research*, pages 457–479, 2004.
- [4] John R Firth. A synopsis of linguistic theory, 1930-1955. *Studies in linguistic analysis*, 1957.
- [5] E Garcia. "a tutorial on okapi bm25". 2011.
- [6] Marti A Hearst. Texttiling: Segmenting text into multi-paragraph subtopic passages. *Computational linguistics*, 23(1):33–64, 1997.
- [7] Karen Sparck Jones. "a statistical interpretation of termspecificity and its application in retrieval". *Journal of Documentation*, 28:11–21, 1972.
- [8] Armand Joulin, Edouard Grave, Piotr Bojanowski, and Tomas Mikolov. Bag of tricks for efficient text classification. *arXiv preprint arXiv:1607.01759*, 2016.

- [9] Taku Kudo. Mecab: Yet another part-of-speech and morphological analyzer. <http://mecab.sourceforge.net/>, 2005.
- [10] Omer Levy and Yoav Goldberg. Neural word embedding as implicit matrix factorization. In *Advances in neural information processing systems*, pages 2177–2185, 2014.
- [11] Quanzhi Li, Sameena Shah, Xiaomo Liu, Armineh Nourbakhsh, and Rui Fang. Tweetsift: Tweet topic classification based on entity knowledge base and topic enhanced word embedding. In *Proceedings of the 25th ACM International on Conference on Information and Knowledge Management*, pages 2429–2432. ACM, 2016.
- [12] Hans Peter. Luhn. "a statistical approach to mechanized encoding and searching of literary information.". *IBM Journal of research and development*, pages 309–317, 1957.
- [13] Tomas Mikolov, Kai Chen, Greg Corrado, and Jeffrey Dean. Efficient estimation of word representations in vector space. *arXiv preprint arXiv:1301.3781*, 2013.
- [14] Lawrence Page, Sergey Brin, Rajeev Motwani, and Terry Winograd. The pagerank citation ranking: Bringing order to the web. Technical report, Stanford InfoLab, 1999.

- [15] 田中 恵 伊藤 孝紀, 深町 駿平, 伊藤 孝行, and 秀島 栄三. ”ファシリテータに着目した合意形成支援システムの検証と評価”. デザイン学研究, pages 67–76, 2015.
- [16] 伊藤孝行 and et al. ”多人数ワークショップのための意見集約支援システム collagree の試作と評価実験”. 日本経営工学会論文誌, pages 83–108, 2015.
- [17] 小谷哲郎, 関一也, 松居辰則, and 岡本敏雄. 好意的発言影響度を取り入れた議論支援システムの開発. 人工知能学会論文誌, 19(2):95–104, 2004.
- [18] 中野哲寛 and 荒木雅弘. 雑談対話システムにおける単語分散表現を用いた話題展開手法. 言語処理学会第 21 回年次大会発表論文集, pages 269–272, 2015.
- [19] 別所克人 et al. 単語の概念ベクトルを用いたテキストセグメンテーション. 情報処理学会論文誌, 42(11):2650–2662, 2001.

付 録 A

発表(予定)論文一覧

A.1 発表(予定)論文一覧

A.2 投稿論文一覧

付 録 B

第 2 回市民共創知研究会

投稿論文

第 2 回市民共創知研究会に投稿した論文を示す。本論文は平成 29 年 7 月 1 日に発表された。

分散表現を用いた話題変化判定

A Topic Change Judgment Method based on Distributed Representation

芳野魁¹ 伊藤孝行¹

Kai Yoshino¹ and Takayuki Ito¹

¹名古屋工業大学情報工学科

¹ Nagoya Institute of Technology, Department of Computer Science

Abstract: As discussions on the Web become bigger, it is expected that discussion of multiple people will be necessary and large scale, and the burden on facilitators will increase accordingly. Therefore, in this research we aim to reduce the burden on facilitators by detecting variance of topics under discussion using distributed representation as one of the burden reductions.

1. はじめに

近年, Web 上での大規模な議論活動が活発になり, 大規模な人数での議論が期待されている. 大規模な議論では意見を共有することは可能であるが, 議論を整理させることや収束させることは難しい. 以上から大規模意見集約システム COLLAGREE が開発された[1]. 本システムでは Web 上で適切に大規模な議論を行うことができるように議論をマネジメントするファシリテーターを導入した.

過去の実験ではファシリテーターの存在が議論の集約に大きな役割を果たしていることが認識されており, 大規模な議論のためにファシリテーターは必要である[2][3]. しかし, 議論の規模に伴って議論時間が長くなる傾向があり, 同時にファシリテーターは常に議論の動向を見続ける必要がある. 故に, 議論の規模が大きくなればなるほどファシリテーターは長時間かつ大規模な議論の動向の監視によって大きな負担がかかる. 大規模な議論が増加する傾向を踏まえるとファシリテーターにかかる負担を軽減する支援が必要となることは明白である.

また, 近年自然言語処理の分野において分散表現が多くの研究で使われており, 機械翻訳を始めとする複数の分野で精度の向上が確認されている[4]. まだ適応されていない分野でも結果の向上が期待できる.

従って, 本研究では負担軽減の1つとして分散表現を用いて議論中での話題の変化を人間の代わりに検知することでファシリテーターの負担を軽減することを目指す.

以下に, 本論文の構成を示す. 第2章では分散表現を用いた話題変化判定を示す. 第3章では評価実

験を行い, 第4章で本論文のまとめを示す.

2. 分散表現を用いた話題変化判定

COLLAGREE を始めとする議論掲示板では, 1つのテーマに対して関連のある複数のテーマを扱う発言が投稿され, 場合によってはある投稿者の発言が親意見となり, 他のユーザーが子意見として返信し, 更に孫意見が存在する.

上記のような議論掲示板での発言に対して, 本文文ではある発言 A と A の子意見, または発言 A と A 直後の発言の間の類似度を計算し, 話題が変化したかの判定を行う手法を提案する. 処理の流れは以下の通りである.

1. 文章の分解
2. 重要度の計算
3. 単語の重み付け
4. 分散表現
5. 類似度の計算

2.1. 文章の分解

文章から単語へ分割するにあたっては形態素解析エンジン mecab を使用した.

2.2. 重要度の計算

「そうですね」や「はい」のような短く, 名詞などの少ない文章は大きな意味は無いが, 他の文章との差異が大きくなってしまいう傾向があったことから文章中の動詞や名詞の数を集計し, 各品詞等の数を基に文章の重要度を求めることを考案した. 提案手

法では重要度

$$\text{Imp}(s) = \frac{\sum_{i \in \text{Pos}} b_i \cdot \text{num}(s, i)}{a \cdot N}$$

$\text{Imp}(s)$: 文章 s の重要度

N : 全文書数

Pos : 品詞集合

$\text{num}(s, i)$: 文章 s 中に現れた品詞 i の数

a : 係数

b_i : 品詞 i に対する係数

を求め、値が閾値 m を下回ったものは重み付けの前に除外する。

係数 a に関しては $a = 2.0$ とし、係数 b_i に関しては

$b_{\text{名詞}} = 1.0, b_{\text{固有名词}} = 2.0, b_{\text{動詞}} = 1.0, b_{\text{形容詞}} =$

$1.0, b_{\text{副詞}} = 1.0, b_{\text{その他}} = 0.5$

閾値 m に関しては $m = 0.45$ とした。

2.3. 単語の重み付け

複数の文書が存在する時、それぞれの文書の特徴付ける単語が特定しにくくなることもある。単語の特定の基準の1つとして TF-IDF という値が使われる。

始めに、TF について説明する。TF は Term Frequency の略で、それぞれの単語の文書内での出現頻度を表し、多く出てくる単語ほど値が大きくなり、重要性が高いことになる[6]。

次に、IDF について説明する。IDF は Inverse Document Frequency の略で、それぞれの単語がいくつかの文書内で共通して使われているかを表す。いくつかの文書で横断的に使われている単語は値が小さくなり、重要性が高くないことになる[7]。

TF-IDF は TF と IDF を掛けたもので、TF-IDF が大きいほどそれぞれの文書の特徴付ける単語であると言える。提案手法では分割された単語の集合の中から TF-IDF が高いものを取り出している。

2.4. 分散表現

自然言語処理において単語の意味を機械に認識させる時、幾つかの方法がある。認識させる方法の1つに単語ごとに人手で意味を付ける方法があるが、人手による手法には幾つか問題点がある。

1. 主観的である。
2. 人間への負担が大きい。
3. 単語間の類似度計算が困難である。

上記の問題を解決するための手法として、単語の言語学的な意味ではなく、文書集合中で周囲に出現している単語の分布を求め、分布を圧縮して密にすることによって単語を低次元の実数値ベクトルで表す方法が考案された[5]。具体的な例を図1で示す。単語の分布に基づき分散表現を使用することで客観的かつ機械の手による意味が付属され、実数値ベクトルであることから数学的な処理が可能になり、単語間の類似度を計算することが可能となっている。分散表現を示した図を以下の図1に示す。

50~300次元				
家	0.9	0.1	0.3	...
犬	0.2	0.8	0.4	...
猫	0.1	0.7	0.4	...
...

図 1: 単語の分散表現

提案手法では分散表現を獲得するにあたり fastText[5] と呼ばれる分散表現への変換手法を使う。学習の際のコーパスには wikipedia の記事データを使用した。

2.5. 類似度の計算

類似度計算においては Python 用の自然言語処理ライブラリ Gensim を使用している。Gensim には単語の集合の類似度を計算する関数を実装されているが本研究で新しい手法を提案する。

類似度を計算する際に単語を集合で比べるよりも一対比較を行った方が精度が高くなると判断したことから提案手法では2つの単語の集合同士をまとめて比較するのではなく、2つの集合中の単語全てで一対比較を行って類似度

$\text{sim}(w1_i, w2_j)$: 単語 $w1_i, w2_j$ 間の類似度

を求め、その平均を文章同士の類似度

$$\text{similarity}(s_1, s_2) = \frac{\sum_{w1_i \in s_1} \sum_{w2_j \in s_2} \text{sim}(w1_i, w2_j)}{\text{num}(s_1) \cdot \text{num}(s_2)}$$

$\text{similarity}(s_1, s_2)$: 単語集合 s_1, s_2 間の類似度

$\text{num}(s)$: 単語集合 s 中の単語数

とする手法を考案した。

3. 実験

3.1. 実験概要

話題変化の検出の実験にあたり、COLLAGREE[3]で取られた「外国人旅行者向けの日本旅行プランに関する議論」の議論データを使用した。実験においては2つのデータと基準を使用した。データ1は議論の進行を支援するファシリテーターが発言しない議論のデータで、データ2ファシリテーターが積極的に発言する議論のデータである。

個々の発言データに話題が変わったかのタグ付けを行い、とある発言AとAに対する返信の発言、またはとある発言AとA直後の発言のどちらかのペアで提案した手法による比較を行い、類似度が閾値を下回った場合を変化ありと検出し、検出された発言が適切かの判別を行った。

評価の際の指針として下記の3つのものを設けた。

1. 正解率:話題が変化したと判断された発言の内、何%が正しく検知されていたか。
2. 網羅率:話題が変わったとタグ付けされた発言の内、何%を正しく検知できたか。
3. 総検知率:全ての発言の内、何%を話題が変化したと検知したか。

評価基準として以下の通り、AとBを設けた。

- A) ファシリテーターの発言は基本的に話題を変える発言であることが多いので検出した発言がファシリテーターに関するものか否かの判別を行い、正誤率や網羅率を評価する。よって、2つの発言の主の一方がファシリテーターであれば正解とする。
- B) 重要であると思われる発言に対してタグ付けを行い、検出した発言にタグがついていた場合正解とする。

また、手法としては下記の3つを使用し、比較した。

- 手法1: 2.5での類似度計算のみを使用する。
- 手法2: 手法1に2.2での重要性推定を追加し、重要でないと思われる単語は事前に除外してから類似度計算を行った。
- 手法3: 手法2に2.3でのTF-IDFによる単語の重み付けを追加し、2つの発言からそれぞれ重要性が高いと思われる単語を抽出してから、類似度計算を行った。

3.2. 実験結果と考察

実験の結果を表1に示す。

表 1: 実験結果

	手法 1	手法 2	手法 3
評価基準 A (データ 1)	正:35.3% 網:19.2% 総:10.7%	正:35.1% 網:7.6% 総:3.3%	正:43.9% 網:22.9% 総:10.4%
評価基準 A (データ 2)	正:63.1% 網:23.5% 総:4.92%	正:64.6% 網:11.3% 総:4.9%	正:71% 網:38.5% 総:15.7%
評価基準 B (データ 1)	正:33% 網:19.2% 総:10.7%	正:50% 網:12.1% 総:3.3%	正:44.5% 網:46.3% 総:10.4%
評価基準 B (データ 2)	正:33.2% 網:31.7% 総:15.4%	正:66.7% 網:23.5% 総:4.9%	正:47.4% 網:60.1% 総:15.7%

手法1から機能を追加していくことで、最終的に多くの場合で総検知率を大きく上げずに正解率、網羅率を上昇させることに成功した。

また、今後の展望としては発言からの単語抽出において更なる工夫が精度を上げるために必要である。

4. まとめ

本研究では分散表現を用いて議論中の話題変化の判定を行った。評価実験により精度が上がってきていることを示した。一方で精度を上げる余地がまだあることも確認した。今後の展望として発言からの単語抽出の改良について現状のTF-IDFを使用した単語抽出ではなく、生成的要約による手法を検証する。

参考文献

- [1] Takayuki Ito, Yuma Imi, Eizo Hideshima, "COLLAGREE: A Facilitator-mediated, Large-scale Consensus Support System", Collective Intelligence, 2015
- [2] 伊藤 孝紀, 深町 駿平, 田中 恵, 伊藤 孝行, 秀島 栄三, ファシリテータに着目した合意形成支援システムの検証と評価, デザイン学研究, 62, 2015, 4_67-4_76
- [3] 伊美裕麻, 伊藤孝行, 伊藤孝紀, 秀島栄三. 大規模意見集約システム COLLAGREE の開発と名古屋市次期総合計画に関する社会実験. 人工知能学会全国大会論文集, 28, 2014, 1-4
- [4] Tomas Mikolov, Quoc V. Le, Ilya Sutskever, "Exploiting

- Similarities among Languages for Machine Translation*",
CoRR, abs/1309.4168, 2013
- [5] Piotr Bojanowski , Edouard Grave , Armand Joulin
and Tomas Mikolov *Enriching Word Vectors with Subword
Information*, 2016.
- [6] Hans Peter Luhn, "A Statistical Approach to Mechanized
Encoding and Searching of Literary Information", *ournal
of research and development. IBM*, 1, 1957, 315
- [7] Karen Sparck Jones, "A Statistical Interpretation of Term
Specificity and Its Application in Retrieval", 28, 1972,
11-21