

単語の共起グラフを用いた潜在的意味に基づく効果的な文書分類 の検証

Validation on Efficient Text Classification Based on Latent Semantic with a Graph of Co-occurring Terms

小倉由佳里* 小林一郎
Yukari Ogura Ichiro Kobayashi

お茶の水女子大学大学院人間文化創成科学研究科理学専攻
Advanced Sciences, Graduate School of Humanities and Sciences, Ochanomizu University

Abstract: We have proposed a method to raise the accuracy of text classification based on latent topic information, introducing several techniques such as extracting important words with PageRank algorithm and reducing the size of target documents by replacing them with important sentences in themselves. We have experimented on text classification with Reuters-21578 data set and confirmed that our proposed method worked to raise the accuracy of text classification. In this paper, we aim to verify our method with additional experiments using 20 Newsgroups data set and report the experimental result.

1 はじめに

近年、インターネットの発達に伴い、爆発的に増大した莫大な量のテキストデータを扱う問題がある。そのため大量のテキストを、自動でカテゴリごとに分類できるような文書分類手法が必要とされている。本研究の先行研究となる [13] では、文書の潜在的意味を考慮した分類手法が提案された。ここでは、文書分類の方針として、まず語彙の重要度に基づき重要文抽出を行い、元の文書を重要文のみで構成し、分類対象となる文書の精錬化を図る。語彙の重要度を定める指標としては、一般に *tf-idf* や語彙の頻度などが用いられるが、語の共起関係からグラフを構成し、PageRank アルゴリズムを用いて重要語の決定が行われた。次に、潜在的意味解析手法を用いて、文書の潜在トピックごとの確率分布をもとに、k-means 法でクラスタリングが行われている。また、実験では Reuters-21578 のデータセットを使用し、提案する手法の有効性を検証した。本稿では、提案する手法の汎用性を検証するために、20 Newsgroups のデータセットを用いた実験結果について報告し、考察を行う。

2 関連研究

文書分類の研究において、分類精度を上げるため数多くの研究がなされており、特に、文書中の語の重要度を定めるアルゴリズムを改良することにより、分類精度の向上が出来ることが報告されている。Hassan ら [1] は、*n*-グラムを用いて、単語の共起関係に基づき構成したグラフに、PageRank アルゴリズムを用いることにより、文書分類の精度が向上することを示した。Zaiane ら [2] や、Wang ら [3] は、文書分類における、語の重要度の決定手法を提案した。Wang ら [3] は、語の重要度の決定に PageRank アルゴリズムを用いることが、文書分類に有効であることを示した。PageRank アルゴリズムは、センチメント分析や、トピック推定にも用いられており、Kubek ら [4] は、語の共起関係に基づき構成したグラフに、PageRank アルゴリズムを用いることにより、トピック推定を行っている。語の重みづけは、文書要約やにおいても重要な課題である。Erkan ら [5] は、LexRank や TextRank と呼ばれる、PageRank アルゴリズムを用いた文書要約の手法を提案している。文をノードとしてグラフを構成し、高い PageRank スコアを持つ、中心性の高い文を抽出することにより、文書要約を行っている。

本研究の先行研究 [13] では、文書を潜在情報に基づいて分類することを目的とし、Newman ら [8] による潜在的情報の首尾一貫性は単語の共起関係により形成されるという報告を参考に、共起語からなるグラフを構

*連絡先：お茶の水女子大学大学院人間文化創成科学研究科理学専攻情報科学コース小林研究室
〒112-8610 東京都文京区大塚 2-1-1
E-mail: g0920509@is.ocha.ac.jp

築し、それに PageRank アルゴリズムを適用することにより、抽出された重要語から重要文を決定する。その重要文を用いて、潜在情報に敏感な文書群を再構成し、文書分類を行う手法を提案した。以下、先行研究 [13] での説明を重複するが、提案手法の内容を再掲しつつ、追加実験の結果と報告と考察を行う。

3 提案方法

3.1 PageRank アルゴリズムによる重要語の決定

PageRank とは、Brin ら [6] によって提案された、Web ページ間に存在するハイパーリンク関係を利用することでページの順位付けを行うアルゴリズムである。PageRank の基本的な考え方は、推薦である。例として、図 1 の場合、 V_a から V_b へリンクが張られているため、これは V_a から V_b への推薦と考えることができる。他の重要な Web ページから推薦されている Web ページは重要である、という考え方が PageRank において中心となっている概念である。Web ページをノード、ページ間のリンク関係をエッジとした有向グラフとして構成され、このグラフに基づいて順位のスコアが計算される。グラフ $G = (V, E)$ が与えられたときに、 $In(V_a)$ は、点 V_a を指している点の集合、 $Out(V_a)$ は、点 V_a が指している点の集合である。点 V_a の PageRank スコアは、式 (1) を反復的に処理することにより、全てのノードの PageRank スコアを求める。 d は、制動係数 (dumping factor) であり、ある一定の割合でリンクのないノードからの影響を考慮するパラメータであり、 $[0, 1]$ の値をとる。

$$S(V_a) = \frac{(1-d)}{N} + d * \sum_{V_b \in In(V_a)} \frac{S(V_b)}{|Out(V_b)|} \quad (1)$$

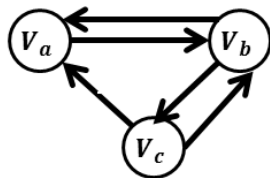


図 1: リンク関係の例

反復計算には、べき乗法を用いる。べき乗法とは、行列の主固有値と主固有ベクトルを見つけるための反復法であり、マルコフ連鎖の定常ベクトルがマルコフ行列の左側主固有ベクトルであること、および、求めた

い PageRank ベクトルが Web ページ間のリンク関係を表した推移行列をもつマルコフ連鎖の定常ベクトルであることにより、PageRank の計算に用いられる。

語の重要度を決定するには、 $tf \cdot idf$ などが頻繁に用いられるが、語同士の様々な関係をグラフ構造で表現し、語の重要度を決定する手法が提案されている [3][1][10]。特に、Hassan ら [1] は、PageRank を用いてランクづけされた語の重要度は $tf \cdot idf$ よりも重要度を明確に差別化できることを示している。本研究でも彼らの手法を参考にして、語の重要度を PageRank アルゴリズムを用いて決定する。

3.2 潜在情報による分類

文書内の潜在的トピックの確率分布を表わすモデルとして Latent Dirichlet Allocation (LDA) [7] がある。このモデルでは、文書内にはいくつものトピックが潜在しており、トピックごとに出現しやすい単語があると考える。各トピックはそのトピックに対する出現確率を持った単語群で表され、複数文書内に存在している総単語に対して、各トピックごとに総和が 1 になる出現確率が割り当てられる。トピック自身にも文書セット内において出現確率の総和が 1 となるトピック比率として確率が付与される。本研究においては、文書に対する潜在トピックの確率分布を用いて、各文書をトピックで構成されるベクトルで表現し、文書間の類似度を測る。

3.3 提案手法における処理の流れ

本手法における、文書分類の流れを説明する。

step1 単語の共起関係の抽出

文書を文で区切り、文脈を考慮して、文中の単語の共起度を自己相互情報量 (PMI: Point-wise Mutual Information) に基づき算出する。

step2 重要単語の決定

step1 で得られた共起関係に基づき、ノードを単語、エッジの重みには PMI を用いたグラフを構成する。図 2 は、共起関係を基に構成したグラフの一例である。ここで、グラフを単語間の PMI で構成する理由は、文書分類を潜在的意味に基づき行うとしており、潜在トピックの一貫性は語の共起関係が影響を与えているとする Newman ら [8] の研究に基づき、潜在トピックを考慮した単語の重要度を算出するためである。このグラフに対し、多くの単語と高い共起度を持つ単語は重要であると考え、PageRank アルゴリズムを用い、単語の重要度のランク付けを行う。

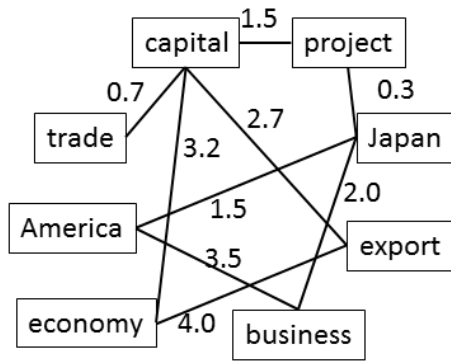


図 2: 類似度グラフ

step3 重要文の抽出

step2 で得られた単語のランキングに基づき、ランキング上位の単語を含む文を重要文とみなし、これを文書から抽出し、元の文書を重要文のみで構成する。

step4 分類

新たに構成された文書群に対し、LDA を用いてそれぞれの文書の潜在トピックごとの確率分布を得る。各文書のトピックに基づくベクトルを Jensen-Shannon 距離を用いて類似度を測り、k-means 法により分類する。

4 実験

4.1 実験仕様

実験対象データには、Reuters-21578¹ のテストデータと 20 Newsgroups² を使用した。提案する手法は、対象文書から重要文を抽出し、文書を精練してから文書分類を行うため、文数の少ない文書では提案手法の効果が判別できないため、1 文書中の文章数が 5 文以上である文書を利用した。

Reuters-21578 のカテゴリは、文書分類の他研究 [9], [11] においても用いられている上位 10 件のカテゴリ、acq, corn, crude, earn, grain, interest, money-fx, ship, trade, wheat を利用した。その結果、文書数 792 件、語彙数 15,835 語、カテゴリ数 10 の文書群を対象に、タグの除去、ステミング処理、ストップワード除去を施し実験を行った。

20 Newsgroups は、20 のニュースカテゴリからなるデータセットである。文書数 11,269 件、語彙数 53,975 語で構成されている。本研究では、文献 [12] を参考にし、

¹ <http://www.daviddlewis.com/resources/testcollections/reuter21578>

² <http://qwone.com/~jason/20Newsgroups/>

comp.graphics, rec.sport.baseball, sci.space, talk.politics.mideast の 4 カテゴリからそれぞれ 200 件ずつをランダムに選び、文書数 800 件、語彙数 14,198 語のデータを使用した。以後、このデータセットを 4-News と記述する。

また、LDA で用いるパラメータは、 $\alpha = 0.5$, $\beta = 0.5$ とし、サンプリングにはギブスサンプリングを用い、イテレーションは 200 回とした。トピック数は、パープレキシティにより決定することにした。トピック数を 1 から 30 まで変化させたときのパープレキシティの値の 10 回の平均をとり、パープレキシティが最小になるときのトピック数を最適トピック数とした。重要文の抽出を行わない元の文書群の分類精度をベースラインとするため、実験に使用するトピック数は最適トピック数を用いた。分類手法には、k-means 法を用い、トピックで構成された文書ベクトルを用いて分類を行う。

4.2 評価手法

評価には、文献 [9] を参考にして、正解率と F 値の 2 つの評価指標を用いる。文書 d_i に関して、 l_i はクラスタリングアルゴリズムにより d_i に与えられたラベル、 α_i は d_i の正解のラベルである。そのとき、正解率は式 (2) で表される。

$$\text{正解率} = \frac{\sum_{i=1}^n \delta(\text{map}(l_i), \alpha_i)}{n} \quad (2)$$

$\delta(x, y)$ は、 $x = y$ ならば 1 となり、そうでなければ 0 となる関数である。 $\text{map}(l_i)$ は、k-means 法により d_i に与えられるラベルである。

評価には、各カテゴリの F 値を求め、全カテゴリの平均を算出した。カテゴリ c_i の F 値は、精度を $P(c_i)$ 、再現率を $R(c_i)$ とすると、式 (3) のように表される。

$$F(c_i) = \frac{2 \cdot P(c_i) \cdot R(c_i)}{P(c_i) + R(c_i)} \quad (3)$$

カテゴリごとの F 値 (式 (3)) を測り、全カテゴリの平均を評価指標として用いた。(式 (4))

$$F = \frac{1}{|C|} \sum_{c_i \in C} F(c_i) \quad (4)$$

また、k-means 法において初期値には、それぞれのカテゴリの正解データの文書ベクトルをランダムに選び、1 つ与えることにする。分類する際、文書群におけるカテゴリ数 k を事前に知っていること、それぞれのカテゴリから 1 つだけ正解例を見つけることは、計算コストがかからないことから、妥当な方法であると判断できる。この方法により、分類結果のクラスが、どのカテゴリであるか判断できるようになる。

4.3 実験結果

k-means 法を 10 回行い、その平均値を測った。ただし、LDA を用いて、文書のトピックごとの確率分布から分類を行う場合には、出力される確率分布 θ が毎回変化する。そのため 1 つの θ に対して k-means 法を 10 回行い、マクロ平均を測った。文書間の類似度指標には、Jensen-Shanon 距離を用いた。重要度の高い上位 3 単語を含む文を抽出して分類を行った場合の正解率と F 値の結果をそれぞれ表 1、表 2 に示す。また、重要文抽出を行った後の文書群の単語数の変化を表 3、表 4 に示す。重要文抽出した後の文書群の単語数が、元の文書群の 8 割程度になるよう設定し、分類を行った場合の正解率と F 値の結果をそれぞれ表 5、表 6 に示す。

表 1: 正解率

単語の重要度	Reuters-21578	4-News
PageRank	0.5671	0.6415
<i>tf · idf</i>	0.5500	0.5915
重要文抽出なし	0.5177	0.8563

表 2: F 値

単語の重要度	Reuters-21578	4-News
PageRank	0.4852	0.6321
<i>tf · idf</i>	0.4347	0.5091
重要文抽出なし	0.4262	0.8494

表 3: Reuters-21578 の単語数の変化

手法	1 語	2 語	3 語	4 語	5 語
PageRank	12,268	13,141	13,589	13,738	13,895
<i>tf · idf</i>	13,999	14,573	14,446	14,675	14,688

5 考察

実験結果の表 1、表 2 より、4-News を用いた実験では、Reuters-21578 を用いた実験と同じ結果は得られなかった。Reuters-21578 では、重要文抽出により文書が精練されたことから、文書の特徴を表現するのに必要な文のみが残り、文書のトピックごとの確率分布の差が測りやすくなったのではないかと考えられた。しかし、4-News を用いた実験では、重要文抽出を行わない場合の方が、行う場合よりも精度が高い結果となった。これは、データセットの性質の違いであると考えられる。この原因としては、4-News は元の文書群の単

表 4: 4-News の単語数の変化

手法	10 語	15 語	20 語	25 語	30 語
PageRank	10,731	10,958	11,078	11,171	11,241
<i>tf · idf</i>	11,048	11,441	11,731	11,849	11,937

表 5: 正解率

単語の重要度	Reuters-21578	4-News
PageRank	0.5529	0.8175
<i>tf · idf</i>	0.5499	0.7948
重要文抽出なし	0.5177	0.8494

語数が Reuters-21578 より少ないことから、重要文抽出を行ったことにより、単語数がさらに減り、本来大量の文書の下で行う学習の効果が下がり、LDA の精度が下がったのではないかと考えられる。また、重要度の高い単語上位 3 単語を含む文を抽出した後の文書群の単語数の変化から考察すると、Reuters-21578 では元の文書群の 8 割程度抽出できているのに対し、4-News では 6 割程度しか抽出できておらず、このため精度が大きく下がったと考えられる。重要文抽出に関しては、Reuters-21578、4-News 共に、*tf · idf* を用いた場合に比べ、PageRank を用いて重要文の抽出を行った場合に文書分類の精度の向上が見られた。このことから、文書の 3 文中での単語の共起関係からグラフを構成し、単語の重要度を PageRank アルゴリズムを用いて決定することにより、分類に適した単語の重要度が得られることが検証された。

また表 3、表 4 から、重要文抽出したあとの単語数の比較では、*tf · idf* と比較して、PageRank を用いた場合に、より語彙数、文数が減っていることが分かる。*tf · idf* の場合、特定の文書に多く出現している単語の値が高くなるため、*tf · idf* が高い単語は、その文書中の多くの文に出現している可能性が高い。そのため、*tf · idf* の高い単語を含む文を抽出すると、自然と多くの文を抽出することになるのではないかと考えられる。Reuters-21578 と 4-News での結果を比較してみると、Reuters-21578 において、重要度の高い上位 1 単語を含む文を抽出した後の単語数と、4-News において、重要度の高い上位 10 単語を含む文を抽出した後の単語数の、元の文書群に占める割合がほぼ等しくなっている。これは、Reuters-21578 は 10 カテゴリであるのに対し、4-News は 4 カテゴリであることから、4-News では、同じカテゴリで似た単語の重要度が高くなっているから抽出単語数が少ないのではないかと考えられる。

表 5、表 6 では、抽出後の単語数を元の文書群の 8 割にして実験を行った。結果は、表 1、表 2 と同じ傾向

表 6: F 値

単語の重要度	Reuters-21578	4-News
PageRank	0.4582	0.8116
$tf \cdot idf$	0.4347	0.7948
重要文抽出なし	0.4262	0.8494

が見られた。これらと比較すると、4-Newsでは、抽出後の単語数の割合を増やしたため、重要文抽出する場合において精度の向上が見られたが、重要文抽出をせず分類を行う場合に一番精度が高くなる結果となった。

6 おわりに

本研究では、先行研究 [13] で提案された PageRank を用いた重要語の抽出を行い、それに基づいて重要文を抽出し、潜在的意味によるクラスタリングを行う手法の汎用性を検証するために、20 Newsgroups のデータセットを用いて、追加実験を行った。実験から、Reuters-21578 では提案手法の有効性が確認されたが、20 Newsgroups を用いた実験では、重要文抽出を行わない場合に最も精度が高くなる結果となり、データセットにより結果に違いが見られた。

今後の課題としては、さらに他のデータセットを用いた実験を行うつもりである。また、文書量が LDA の精度に影響することが考えられることから、文書数をさらに増やした実験を行いたいと考えている。さらに、トピック数を変化させた場合の実験結果の比較を行う。また、現在は k-means 法での分類しか行っていないため、他の多クラス分類手法との比較を行うつもりである。

参考文献

[1] Samer Hassan, Rada Mihalcea, Carmen Banea.: Random-Walk Term Weighting for Improved Text Classification, (2007)

[2] Osmar R.Zaiane, Maria-luiza Antonie.: Classifying Text Documents by Associating Terms with Text Categories, *In Proc. of the Thirteenth Australasian Database Conference(ADC'02)*, pp. 215–222

[3] Wei Wang, Diep Bich Do, and Xuemin Lin.: Term Graph Model for Text Classification, *Springer-Verlag Berlin Heidelberg 2005*, pp. 19–30 (2005)

[4] Mario Kubek, Herwig Unger.: Topic Detection Based on the PageRank's Clustering Property, *IICS'11*, pp. 139–148 (2011)

[5] Gunes Erkan.: LexRank: Graph-based Lexical Centrality as Salience in Text Summarization, *Journal of Artificial Intelligence Research 22*, pp. 457–486 (2004)

[6] Sergey Brin, Lawrence Page.: The Anatomy of a Large-scale Hypertextual Web Search Engine, *Computer Networks and ISDN Systems*, pp. 107–117 (1998)

[7] David M. Blei, Andrew Y. Ng, Michael I. Jordan, John Lafferty.: Latent dirichlet allocation, *Journal of Machine Learning Research*, Vol. 3, p. 993–1022 (2003)

[8] Newman David, Lau Jey Han, Grieser karl, Baldwin Timothy.: Automatic evaluation of topic coherence, *Human Language Technologies :The 2010 Annual Conference of the North American Chapter of the Association for Computational Linguistics*, pp. 100–108 (2010)

[9] Gunes Erkan.: Language Model-Based Document Clustering Using Random Walks, *Association for Computational Linguistics*, pp. 479–486 (2006)

[10] Christian Scheible, Hinrich Shutze.: Bootstrapping Sentiment Labels For Unannotated Documents With Polarity PageRank, *Proceedings of the Eight International Conference on Language Resources and Evaluation*, (2012)

[11] Amarnag Subramanya, Jeff Bilmes.: Soft-Supervised Learning for Text Classification, *Proceedings of the 2008 Conference on Empirical Methods in Natural Language Processing*, pp. 1090–1099, Honolulu (2008)

[12] Liping Jing, Michael K.Ng, Jun Xu, Joshua Zhexue Huang.: Subspace Clustering of Text Documents with Feature Weighing K-Means Algorithm *PAKDD 2005, LNAI 3518*, pp. 802–812 (2005)

[13] 小倉 由佳里, 小林 一郎. : 単語の共起グラフを用いた潜在的意味に基づく効果的な文書分類への取り組み, *インタラクティブ情報アクセスと可視化マイニング第3回研究会*,(2013)