

共起と出現確率と単語間距離に基づく連語抽出は、 ユーザプロファイルの自動名づけを可能にする

望月 啓太[†] 清水 康太郎[†] 毛利 瞳[†] 金道敏樹[†]

金沢工業大学工学部情報工学科[†]

1. はじめに

我々は、コンテンツベースの推薦システムに用いられるユーザプロファイルに対する自動名づけ機能を実現することを目指している。

推薦システムは利用者にとって有用と思われる対象、情報、または商品などを選び出し、それらを利用者の目的に合わせた形で掲示するために利用される[1]が、現状ユーザプロファイルに対する自動名づけは、

- ・ユーザの興味の可視化
- ・ユーザプロファイルの相互作用の利便性向上
- ・ユーザプロファイルの興味毎の分割
- ・ユーザの興味の遷移の可視化

につながる有用な機能である(図1)にも関わらず、残念ながら行えていない。従来研究も、「Paragraph Vectorに基づくニュース記事分類のための自動ラベリング」[2]のような限定的な研究があるだけである。



図1 嬉しさのイメージ

今回、複数の単語を組合せ、連語を探す方法で、ユーザプロファイルに対する自動名づけアルゴリズムに取り組んだので報告する。アルゴリズムの開発にあたり、我々は、情報フィルタ INSOP を用いた論文推薦システムを前提とした。INSOPは、軽量コンパクトなアルゴリズムであり、また、プロファイルの分割が可能な点が魅力的である。推薦システムのコンテンツとして「論文」を選択した理由は、論旨がはっきりしている文章から検証を開始したかったことによる。

2. 提案手法

我々が提案する自動名づけアルゴリズムは、連語インデキシングと代表連語抽出からなる(図2)。代表連語抽出は、2つのサブアルゴリズムからなる。第1アルゴリズムは、出現確率を用いるものであり、第2アルゴリズムは、単語間距離を用いるものである(図3)。

2.1. 連語インデキシング

論文におけるユーザプロファイルを名付けるためのレベルに適するようなキーとなる概念を表す連語の抽出アルゴリズムを図4に示す。

Two connected nouns supported by occurrence probability, co-occurrence and relative location in documents enables automatic labeling of user profiles

Mochizuki Keita[†] Shimizu Kotaro[†] Mouri Hitomi[†] Kindo Toshiki[†]
Information and Computer Science, the College of Engineering, Kanazawa Institute of Technology[†]

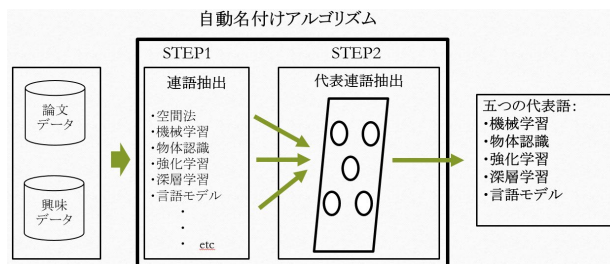


図2 自動名づけアルゴリズムの概略図



図3 第2アルゴリズムの概略図

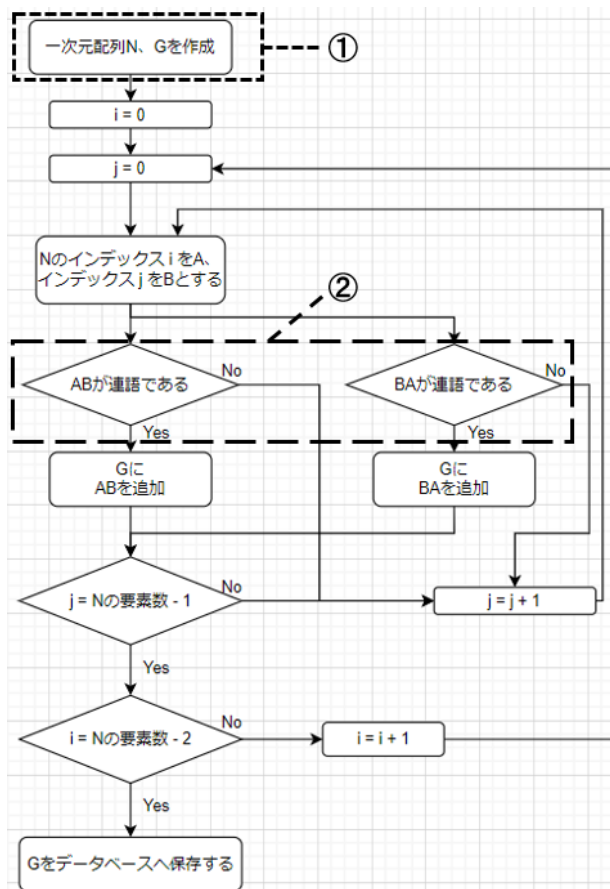


図4 連語抽出アルゴリズムのフローチャート

図4のフローチャートに沿って説明を行う。まず、①内の一次元配列Nは、ある論文の名詞だけで構成された

共起データにおける他の名詞との結びつく数(リンク数)が多い名詞が配列の先頭になるようにデザインされた配列である。また、①内の一次元配列Gは、抽出された連語を保存するための配列である。次に②内で行われる「BA, ABの連語判定」では、Wikipediaを活用している。

図4のアルゴリズムは、対象の論文P(i), i=1, 2, …全に対して適用され、各論文に対して、連語インデックスI(P(i))を作成する。

2.2. 代表連語の抽出

代表語の抽出は2.1.で抽出した論文毎の連語インデックスから、ユーザプロフィールのラベルとなる連語を抽出することである。今回、複数の論文を2つの集合に分割したとき、一方の集合に現れる確率を用いた方式と、論文中の単語の出現位置を用いる方式を組み合わせアルゴリズムを構築した。

確率を用いた方式では、任意に2分割した集合からその分割に必要な連語を抽出できる。

まず、連語インデックスに含まれる連語の集合Uを考える。次に論文集合Sを二つの集合S1とS2に分割する。Uの内、ある連語aがS1に出現する確率を求め、Ps(Probability score): $Ps = \frac{\sum P(S1)|a}{P(U)|a}$ として算出する。Psが1に近ければS1を代表する連語、0に近ければS2を代表する連語とできる(図5)。

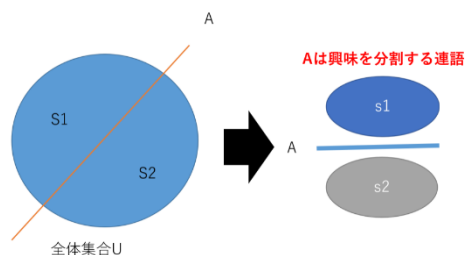


図5 集合Uを2分割する連語の抽出

これにより、あるユーザの興味の有無を含んだユーザプロフィールが与えられたとき、上述のS1にINSOPを用いてユーザが興味を持つもの、S2を、興味を持たないものとするれば、興味の内容に関わらず、ユーザが何に興味があり、興味が無いかを示す連語を可視化できる。

この出現確率を用いたフィルタリング機能だけでは時折、集合S1とS2の関係を表さない連語が残る。その原因は連語の”自然さ”はWikipediaの検索語であることで保証されているものの、共起しか調べていないため論文の内容を反映しているとは限らないからである。ここでは連語をつくる2つの単語が論文中に近い「距離」で現れるかどうかを調べて、距離スコアDs(Distance score)とする。

名詞Aと名詞Bを組み合わせた連語をABとしたとき論文内の名詞A, Bの間に入っている名詞の数を単語間距離dと設定し、dがより短い方がラベリングに妥当な連語であると考えDsを($Ds = \sum(100 - 5d)n, n = \text{出現回数}$)と定義した。

連語の最終的なスコアは、両者の積(score × Ds)となり、その値が大きいもの5つを代表連語とする。

3. アルゴリズムの検証と考察

8人のユーザに自動名づけアルゴリズムに対する満足度を回答して貰い、アンケートによるアルゴリズムの性能評価を行った。

自動名づけアルゴリズムに対するユーザの満足度の評価方法は、まず、8人のユーザに、情報工学に関する300本の論文の内、人工知能に関する論文10本程度に対して興味の有無を回答してもらい、ユーザプロフィールを作成する。次に、ユーザプロフィールに対してアルゴリズムを適用し、アルゴリズムの出力である代表連語(5つの連語)をユーザに見て貰い、その連語がユーザの興味を表現できているのかを10段階評価で評価をしてもらった(図6)。

10段階評価:ユーザの興味を表現する連語が

「1: 全くない ~ 10: 妥当なものがある」

評価の結果、ユーザの評価の平均が8.25点と高く、アルゴリズムを用いて抽出された連語の妥当性が示された。

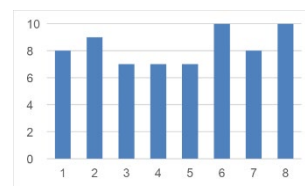


図6 ユーザの連語に対する評価 {x: UserID, y: points}

UserID4と6のユーザプロフィールからは以下の連語が抽出され、人工知能というジャンルの中でも興味の違いが確認でき、「machine」と「learning」それぞれ一語では意味が広く漠然としてしまうが、「machine learning」にすることで意味が判然となり連語の強みも確認できた(表1)。

表1 UserID4と6の代表連語

4	information model	language model	object detection	learning model	model approach
6	reinforcement learning	machine learning	image fusion	user interface	change detection

上記の実験結果から、自動名づけアルゴリズムによるユーザプロフィールのラベリングは、有効であり、ユーザの興味を可視化できると結論できる。

4. まとめ

ユーザの興味の有無を利用した自動名づけは、ユーザの興味をユーザプロフィールにラベリングすることで可視化できる。また、興味の可視化によりユーザプロフィールを利用したシステムの利便性向上のチャンスにつながることを示された。

5. 謝辞

「距離尺度を用いた意味のある連語の選択」部分において学部3年の毛利瞳さんに協力していただきました。心より感謝申し上げます。

参考文献

- [1] Toshiaki KINDO et al. Adaptive Personal Information Filtering System that organizes personal profiles automatically, Proceeding of IJCAI97, 1997.
- [2] 斎藤大志, 内田理. Paragraph Vector に基づくニュース記事分類のための自動ラベリング, 電子情報通信学会総合大会, 2017