

トピックモデルにおける Latent Dirichlet Allocation を用いた 学術論文情報の可視化

宇佐美宏樹^{†1} 高岡詠子^{‡2}
上智大学大学院^{†1} 上智大学^{‡2}

1. 研究背景

現在, 日本の医師を取り巻く環境は問題を抱えており, 特に論文執筆に関してはより重視しながら環境改善に取り組まなければならない. 日本と米国の年代別における論文執筆数を見ると, 各年代で米国よりも執筆数は劣っているが, 特に「研修医」と呼ばれる20年代に大きな差が見られる[1]. この様な現状の原因として, 研修医が1日の中で研究や自己研磨に当てることができる時間が非常に少ないこと[2]が挙げられる. そこで, 現在の環境の中で, 短時間で効率的に論文執筆の際に行う先行研究の調査や情報収集ができる様なソリューションが求められている.

2. 研究目的

上記の背景を受けて本研究の目的を, 「医学論文の作成を行う日本の若い医師達が, 短時間で効率よく資料を検索・収集できる様な環境を提供すること」とする.

3. 研究内容

多くの医師は, PubMed と呼ばれる無料論文検索エンジンを使用し, 先行研究や実験結果の調査を行っている. 研究や自己研磨の時間を短くするにはどうすれば良いかという点で, 着目したのがこの資料検索に当てる時間である.

特に, 学術論文の探索的検索において初期段階で行う文献レビューは, 膨大な論文数の中から, 希望する論文にたどり着くまで, 何度もレビューを繰り返さなければならず, 時間的コストが高い. 特に若い医師や研究者, 未知の研究分野に取り組む研究者にはこの文献レビューはより負担となる.

3.1 関連研究

Claus Boye Asmussen 氏らは, 機械学習の手法の進歩によって, 初期段階の文献レビューの時間的コストを抑えながら, 多くの問題に対処することができるとした[3]. 取り入れた機械学習の手法は, 教師なし手法の分野であるトピックモデリングの1つ LDA (Latent Dirichlet Allocation) であり, R を用いて探索的な文献レビューのあらゆる文脈に適用でき, 完全な文献レビューにも使用できる汎用的なフレームワークを作成した[4].

また Zineb Sabri 氏らは, LDA と Python を用いた教師なし

トピックモデリング法により, 複数の科学データベースから400以上の研究論文を長年にわたって収集し, 処理し3つの主要トピックを特定した[4]. 先行研究を参考にから, Pubmed における abstract 情報をデータセットとし, LDA と python パッケージである pyLDAvis を使用して, 論文情報の可視化に取り組むこととした.

3.2 研究手法

今回の論文情報の可視化においては以下の手順で行った.

- ① NCBI (National Center for Biotechnology Information) のデータベースに対するユーザ向けの取得システム Entrez と Biopython を使用し, Pubmed 上にある論文の abstract 情報の取得
- ② 取得した abstract 情報から, ストップワード及び, 記号, 数字の除去
- ③ 全文章において, 指定した出現回数に満たない単語の除去
- ④ LDA モデリングに必要な辞書とコーパスの準備. 辞書に関して, Bag-of-Words を用いてベクトル表現に変換し, さらに, 特定の文書にしか出現しない単語と文書全体で使用頻度が低い単語の除去を行った.
- ⑤ トピックス数の検討とモデリング
モデルの評価指標として一般的に用いられる Perplex と Coherence を使用しトピック数の検討を行い, Perplexity の値と Coherence の値の差が大きくなることをモデルのトピック数とした.
- ⑥ 選択したトピック数で LDA モデルを行い, pyLDAvis を使用して可視化を行った. 今回の可視化を行うにあたり, Metric Multi-dimensional Scaling (MMDS)の尺度を利用して可視化を行った.

今回1検索クエリあたりの論文情報の取得件数を10000件とした. これは Entrez の EFetch におけるレコード取得数の最大値が10000件であるからである. 図1は, 一例として検索クエリを ICI (Immune Checkpoint Inhibitor) とし10000件の論文データからトピックの可視化を行ったものである.

3.3 実証実験

作成した論文情報の可視化が, 探索的検索の先行レビューにおいてどの程度軽減できるのかを確かめるために, 以下の実証実験を行った.

^{†1} HIROKI USAMI, Sophia University Graduate School of Science and Technology / Information Science .

^{‡2} EIKO TAKAOKA, Sophia University .

