

エビデンスに基づく医療のための 文献キュレーションシステムの開発

小林 賢司^{1,a)} 阿部 修也^{†1,b)} 森田 一^{1,c)} 吉川 和^{1,d)} 岡嶋 成司^{1,e)} 富士 秀^{1,f)}

概要：本稿では、自然言語処理技術とアノテーションツールを連携させることにより、医療文献のキュレーション作業を効率化するシステムについて述べる。キュレーション作業は、治療法のエビデンスとなる医療文献の中から、疾患名や薬剤名など医療情報を抽出する作業であり、「エビデンスに基づく医療」を実現する上で必要不可欠である。しかし、膨大な医療文献の中から必要な情報を抽出することは容易ではなく、多大な労力を要する。そこで本研究では、自然言語処理技術によって抽出処理を自動化し、その抽出結果をアノテーションツールにより確認・編集を行えるシステムを開発した。

Biomedical Literature Curation System For Evidence-Based Medicine

KENJI KOBAYASHI^{1,a)} SHUYA ABE^{†1,b)} HAJIME MORITA^{1,c)} YOSHIKAWA HIYORI^{1,d)} SEIJI OKAJIMA^{1,e)}
MASARU FUJI^{1,f)}

1. はじめに

医療業界において、「エビデンスに基づく医療」の確立が推進されている。エビデンスとなる文献が十分に存在する治療法であるか踏まえた上で、患者の状況に合わせて治療を行う医療のあり方である。エビデンスには、エビデンスレベルと呼ばれる信頼性の高さを示す指標があり、特にバイアスがなく大規模に行われた研究は、エビデンスレベルが高く評価される。そのようなエビデンスレベルの高い文献の治療法、つまり十分に治療効果や副作用などが証明された治療法を選択することで、患者にとって最も有益で、かつ害の少ない結果となることが期待できる。

しかし実現するには、医療者が日々進歩する医療情報を把握する必要がある。これは膨大な量の文献に対して容易

ではない。そのため、エビデンスレベルの高い文献に絞ることや、中身を参照せずとも文献の内容を把握できる仕組みを提供することが望ましい。これを行うためには、文献から抽出した情報を整理したデータベースの構築が必要である。

このようなデータベースを構築する作業はキュレーションと呼ばれ、また医療専門知識を有するキュレーターと呼ばれる者によって盛んに行われている。キュレーターは優先度の高い文献を選別し、文献を読み解く中で必要な情報を抽出し、十分に吟味した上で抽出した情報をデータベースに登録する。この作業は、医療文献が日々進歩・増大することも相まって、多大な人手と期間を要する。

一方、近年 AI 技術は目覚ましく進歩しており、様々な分野において適用効果が見込まれている。この医療文献のキュレーションにおいても、大量の文章を対象とすることから、AI 技術の一つである自然言語処理技術を適用すれば、自動抽出によってキュレーション作業負荷を低減することが期待できる。

しかし、自然言語処理技術で機械的に抽出した情報をそのまま医療現場に適用することは適切ではない。また、精度が低ければ却ってキュレーション作業の効率が下がりが

¹ 株式会社富士通研究所
Fujitsu Laboratories LTD., Kanagawa, 211-8588, Japan
^{†1} 現在、富士通株式会社
Presently with Fujitsu Limited, Kanagawa, 211-8588, Japan
a) kobayashi.kenji@fujitsu.com
b) abe.shuya@fujitsu.com
c) hmorita@fujitsu.com
d) y.hiyori@fujitsu.com
e) okajima.seiji@fujitsu.com
f) fuji.masaru@fujitsu.com



図 1 マニュアルキュレーション作業手順

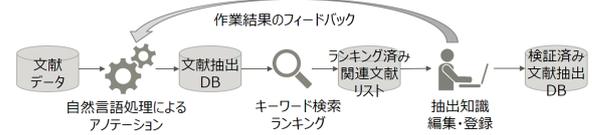


図 2 提案システム作業手順

ねない。医療者同様に自然言語処理技術としても、医療分野の新しいトピックやトレンドを追従していかなければ、精度はさらに下がっていくであろう。

そこで本研究では、自然言語処理技術による自動抽出を行った上で、キュレーターが抽出情報の確認や修正するためのインターフェースを提供し、さらに再学習用に修正内容を学習データとしてフィードバックする文献キュレーションシステムを開発する。これにより、キュレーターの作業を低減した上で、現場への知識適用、および継続的かつ高精度な抽出を実現する。

2. 医療文献キュレーション

本章では、キュレーションの目的や作業内容について述べる。また、従来のキュレーション作業（マニュアルキュレーション）の課題と、提案システムが想定する自然言語処理技術を適用した場合のキュレーション作業について述べる。

2.1 医療文献キュレーションの目的

本研究で扱うキュレーションシステムは、遺伝子疾患に関連する文献を対象としており、特に、遺伝子変異と疾患に関する文献中の知識を整理することを目的としている。キュレーション作業において、抽出・整理の対象とする知識は、以下の通りである。

- (1) 疾患名、薬剤名、遺伝子名、遺伝子変異名の 4 タイプの固有表現
- (2) 固有表現間の関係情報
- (3) 実験手法、分析手法

(2) の固有表現間の関係とは、例えば、薬剤 A は疾患 A に「有効」である、といった抽出された固有表現の臨床的意義に関わる情報であり、(3) は、何を対象とした実験か、どういった実験結果の分析を行っているか、など、内容の信頼性に関する情報である。キュレーションによって蓄積された情報は、クリニカルシーケンスにおいて、治療方針やゲノムレポートを作成する際の参考情報として利用したり、創薬においてターゲットとする遺伝子変異を探索するために活用するなど、様々な応用が考えられる。

マニュアルキュレーションでは、キュレーターが対象とする遺伝子名や疾患名を利用して文献の検索を行い、アブストラクトを利用して対象文献を絞り込んでから、実際に文献を読み、遺伝子変異の病原性や、疾患に対する薬剤の有効性といった知識を抽出する（図 1）。文献に記載された知識の中には、新しい知識としてデータベースに登録する

にはエビデンスが不十分なものも多く存在する。そのため、キュレーターは、文献から知識を網羅的に抽出するだけでなく、エビデンスレベルに関する指標などを参照しながら、内容の信頼性を検証し、知識の扱いを判断する必要がある。作業負荷が非常に大きい。また、遺伝子名や疾患名には多くの表記ゆれや語義の曖昧性があり、関係情報の取りこぼしや誤りが発生し得る。

提案システムは、自然言語処理技術を適用することで、事前に各文献に含まれる固有表現や固有表現間の関係情報にアノテーションを行っておき、キュレーターが文献を選択した際に、文献に含まれる固有表現やその関係情報を周辺の記述とともに提示する。キュレーターは提示された情報を確認し、編集・登録作業を行う。各文献にはエビデンスレベルに基づいたスコアが付与されており、キュレーターがキュレーション対象とする遺伝子名や疾患名を利用して文献の検索を行うと、スコア順に文献が提示される。そのため、キュレーターは信頼性の高い文献から優先的にキュレーション作業を進めることができる（図 2）。さらに、疾患名や遺伝子名といった固有表現は、抽出と同時に既存の医療系知識ベースの ID と関連付けられており、文献検索時のキーワードも ID に変換して検索に利用されるため、表記ゆれや曖昧性に頑健な検索を行うことができる。また、キュレーターの登録した情報は、自然言語処理のモデル生成に利用するための学習データに変換され、定期的にモデルの再学習が行われる。これにより、継続的なアノテーション精度の向上と、新語・未知語への対応を図る。

3. 関連研究

キュレーションの支援ツールとして、アノテーションツールが挙げられる。[1] では、30 種の医療文献向けのアノテーションツールを挙げ、特にその内 13 種のツールについては、ドキュメント、データ形式、機能、性能といった観点について 35 項目で評価している。しかし、いずれのツールについてもエビデンスレベルを判断する機能をサポートしていない。一方、エビデンスレベルを判断する研究として、[2], [3], [4], [5], [6] が挙げられる。これらは逆に、アノテーションツールとの連携は考慮されていない。従って、本研究では、エビデンスに基づく検索と自然言語処理を利用したアノテーション支援システムをシームレスに連携させ、信頼性の高い知識の蓄積を目指す。

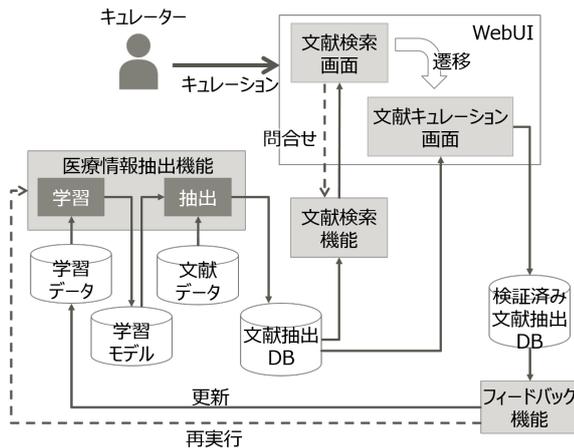


図 3 文献キュレーションシステム概要図

4. 文献キュレーションシステム

ここでは、文献キュレーションシステム全体の概要について述べた後、各画面 UI や機能について詳しく述べる。

4.1 システム概要

文献キュレーションシステムの概要図を図 3 に示す。文献キュレーションシステムは、フロントエンドにキュレーターが作業するための WebUI を備え、バックエンドで文献抽出 DB を構築する。

フロントエンドの WebUI は、文献検索画面と文献キュレーション画面から成る。文献検索画面は、キュレーターがキュレーションしたい文献を検索するための画面である。固有表現や関係を指定して検索すると、文献の一覧が表示され、各文献のタイトルやアブストラクトの他、自然言語処理により抽出した関係や分類などが表示される。キュレーターはそれらの情報を参考にして文献を選択すると、文献キュレーション画面に遷移する。文献キュレーション画面では、アノテーションされた文献の文章が表示される。そこでキュレーターは、誤りがあれば訂正し、確認を終えたら修正内容を保存する。

バックエンドでは、自然言語処理やキュレーション作業によって、文献抽出 DB の作成や更新を行う。文献抽出 DB は、固有表現や関係、分類など自然言語処理によって文献から抽出された情報が格納されたデータベースである。この文献抽出 DB は、医療情報抽出機能によって構築される。医療情報抽出機能は、アノテーションされた学習データから学習モデルを構築し、それを使用してアノテーションされていない文献データ、つまり本システムにおいてキュレーション対象となる文献データから医療情報を抽出し、文献抽出 DB に格納する。文献検索機能では、文献抽出 DB を使用して、検索リクエストに関連する文献を検索結果として返す。また文献キュレーション画面によって

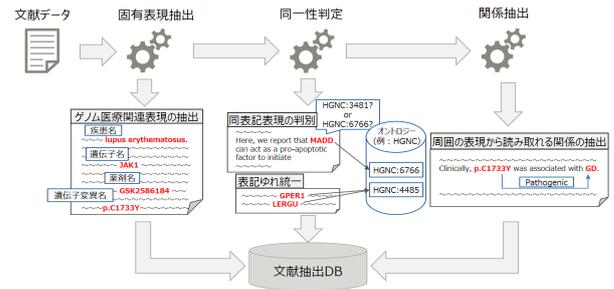


図 4 自然言語処理による医療情報抽出

検証された文献は、フィードバック機能によって再学習される。このように、学習モデルを再構築、再適用することでより精度を高めた抽出を可能とする。

4.2 医療情報抽出機能

医療情報抽出機能は、事前に各文献に対して、固有表現抽出、同一性判定、関係抽出を適用し、テキストに対するアノテーションを行う (図 4)。固有表現抽出は、文献から疾患名 (Disease)、薬剤名 (Drug)、遺伝子名 (Gene-Protein)、遺伝子変異名 (Mutation) を抽出するために適用され、各単語の周囲の表現から、対象とする固有表現であるかどうかを判定する。同一性判定は、抽出された疾患名、薬剤名、遺伝子名を、Disease Ontology[7]、DrugBank[8]、HGNC[9] といった既存のオントロジー/DB の ID に紐づけるために適用され、事前学習した周囲の文脈から、同表記で意味の異なる表現や表記ゆれの問題を解決して、対応する ID を特定する。関係抽出は、文献に含まれる固有表現間の関係を抽出するために適用される。本システムでは「遺伝子と疾患」、「遺伝子と薬剤」、「薬剤と疾患」、「遺伝子変異と疾患」、「遺伝子変異と薬剤」の 5 種類の関係を対象として、周辺文脈から関係の有無を判定する。

4.3 文献検索画面

文献検索画面は、キュレーターが最初に操作する画面である (図 5)。画面上部に備える検索フォームは、一般的な用語だけ指定するのではなく、2 個の固有表現とそのタイプ、それらの関係を指定する。例えば図中のように固有表現が「Disease - Drug」であれば、一つ目の固有表現「tumor」は「疾患名」、二つ目の固有表現「gemcitabine」は「薬剤名」であることを意味する。タイプは他にも「Gene-Protein - Disease」、「Mutation - Drug」など、遺伝子や遺伝子変異についてもプルダウンメニューから選択できる。また関係については、2 個の固有表現間に有効であることを示していれば「Positive」、そうでなければ「Negative」を指定する。従って、図中の例は、「疾患名:tumor」に対し、「薬剤名:gemcitabine」が有効であることを示す文献を検索することを意味する。

画面下部の表は、検索結果となる。左の列から、タ

| 固有表現 | | Disease - Drug | | tumor | gemcitabine |
|--|---|--|-------------|---------|-------------|
| 関係 | | Positive | | | |
| 検索 | | | | | |
| Title | Citation | Relation | Reliability | | |
| Recent advances in radiation therapy of pancreatic cancer | F1000Research F1000 Research Limited 7- -2018-12-13 | GVAX Drug Positive Disease ----> cancer SBRT Drug Positive Disease ----> tumor 5-fluorouracil Drug Positive Disease ----> BRPC FOLFIRINOX Drug Positive Disease ----> BRPC oxaliplatin Drug Positive Disease ----> BRPC | ★ | ★ | ★ |
| Pancreatic cancer has a dismal prognosis with an overall survival outcome of just 5% at five years. However, paralleling our improved understanding of the biology of pancreatic cancer, treatment paradigms have also continued to evolve with newer advances in surgical techniques, chemotherapeutic agents, radiation therapy (RT) techniques, and immunotherapy paradigms | | model(00000000) | Relevance | 解析手法 | 実験対象種 |
| | | | 統計的解析 | 人 | 実験対象レベル |
| | | | 遺伝子解析 | マウス・ラット | 個体/組織 |
| | | | タンパク質解析 | その他の動物 | 培養細胞 |
| | | | 細胞解析 | | タンパク・生体分子 |
| | | | その他 | | |
| Efficacy of the Oral Fluorouracil Pro-drug Capecitabine in Cancer Treatment: a Review | Molecules MDPI 13-8- 2008-8-27 | 5-FU Drug Positive Disease ----> colorectal cancer Disease LV Drug Positive Disease ----> colorectal cancer Disease Capecitabine Drug Positive Disease ----> breast cancer Disease docetaxel Drug Positive Disease ----> breast cancer Disease oxaliplatin Drug Positive Disease ----> colorectal cancer Disease | ★ | ★ | ★ |
| Capecitabine (Xeloda?) was developed as a pro-drug of fluorouracil (FU), with the aim of improving tolerability and intratumor drug concentrations through its tumor-specific conversion to the active drug. The purpose of this paper is to review the available information on capecitabine, focusing on its clinical effectiveness against various carcinomas. Identification of all eligible English trails was made by searching the PubMed and Cochrane databases from 1980 to 2007 | | model(00000000) | Relevance | 解析手法 | 実験対象種 |
| | | | 統計的解析 | 人 | 実験対象レベル |
| | | | 遺伝子解析 | マウス・ラット | 個体/組織 |
| | | | タンパク質解析 | その他の動物 | 培養細胞 |
| | | | 細胞解析 | | タンパク・生体分子 |
| | | | その他 | | |
| Potential of the Dietary Antioxidants Resveratrol and Curcumin in Prevention and Treatment of Hematologic Malignancies | Molecules MDPI 15-10- 2010-10-12 | resveratrol Drug Positive Disease ----> lymphocytic leukemia Disease curcumin Drug Positive Disease ----> lymphoma Disease curcumin Drug Positive Disease ----> T-cell leukemia Disease Curcumin Drug Positive Disease ----> cancers Disease resveratrol Drug Positive Disease ----> leukemic Disease | ★ | ★ | ★ |
| Despite considerable improvements in the tolerance and efficacy of novel chemotherapeutic agents, the mortality of hematological malignancies is still high due to therapy relapse, which is associated with bad prognosis. Dietary polyphenolic compounds are of growing interest as an alternative approach, especially in cancer treatment, as they have been proven to be safe and display strong antioxidant properties. Both polyphenols are currently being tested in clinical trials | | model(00000000) | Relevance | 解析手法 | 実験対象種 |
| | | | 統計的解析 | 人 | 実験対象レベル |
| | | | 遺伝子解析 | マウス・ラット | 個体/組織 |
| | | | タンパク質解析 | その他の動物 | 培養細胞 |
| | | | 細胞解析 | | タンパク・生体分子 |
| | | | その他 | | |

図 5 文献検索画面

イトル・アブストラクト、引用、固有表現・関係、信頼性を示しており、1 文献につき 1 行で表示する。固有表現・関係カラムは、その文献から抽出された固有表現・関係をリスト化している。信頼性カラムは、Relevance が検索クエリとの関連性、その右の解析手法や実験対象種、実験対象レベルは文献が持つ分類であり、それらを加味した総合スコアを星で 5 段階により表現している。

4.4 文献キュレーション画面

文献検索画面からキュレーションしたい文献を選択すると、文献キュレーション画面が表示される (図 6)。アノテーションツールである brat[10] をベースに開発し、固有表現はハイライトされ、関係は矢印で結ばれる。各固有表現上部にはそのタイプが示され、関係の矢印間には関係のタイプが示される。例えば、図中 1 行目の「pancreatic cancer」は「疾患名」、「gemcitabine」は「薬剤名」を表し、関係の矢印により「gemcitabine」が「pancreatic cancer」に「Responsive」であることを意味する。

これらの固有表現と関係について、語句の範囲、固有表現または関係のタイプを編集することが出来る。例えば図 7 は、固有表現を選択した際の編集画面であり、「Entity Type」の項目からタイプを選択できる。また抽出漏れがある場合も、テキストを選択することで編集画面が表示さ

れ、登録処理を行える。固有表現間をドラッグで結ぶことで関係の追加も可能である。

このような確認および編集を一通り行うことで、検証済み文献抽出 DB に保存される。

4.5 文献検索機能

文献検索機能は、文献検索画面の検索リクエストに対し、関連する文献の一覧を結果として返す。文献一覧の結果は、リクエストされた固有表現や関係に関連する文献となる。また、その中でもエビデンスレベルの高い文献を上位にして返している。

エビデンスレベルは、文献で述べられる実験内容の信頼度を表す。一般的な文献の評価指標としてインパクトファクター [11] やオルトメトリクス [12], h-index[13] などが挙げられるが、論文の被引用数をもとに算出する指標であり、文献に記述される治療法の信頼度を表すわけではない。一方、エビデンスレベルは、大規模かつバイアスのない実験方法であるほど高く評価される。代表的な指標としてオックスフォード大学が提唱した指標 [14] が挙げられるが、例えば Randomized Controlled Trial 手法のシステマティックレビューは、バイアスのない実験を行った文献について網羅的に評価している文献であり、エビデンスレベルは最も高い。

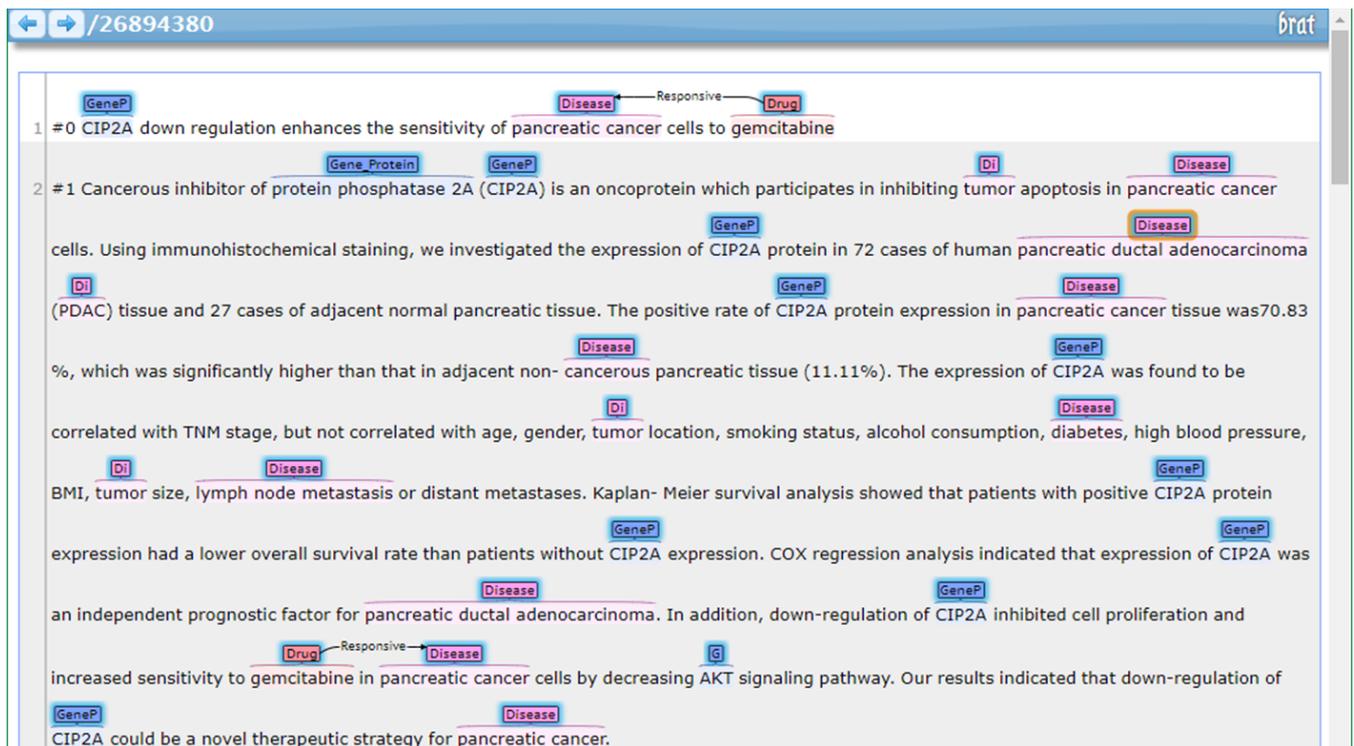


図 6 文献キュレーション画面



図 7 文献キュレーション 編集画面

文献検索機能はこの考えを取り入れて、実験内容の分類をもとにスコアリングしている。

4.6 フィードバック機能

フィードバック機能では、検証済みの文献抽出 DB から、学習データを更新する。キュレーターによって検証された質の高い学習データを追加することが出来るため、再構築した学習モデルでは、より高精度な抽出を行える。ただし、再構築した学習モデルをすぐ適用することには課題

表 1 学習データ内訳

| 処理 | 件数 | 単位 |
|--------|-----------|---------|
| 固有表現抽出 | 15,091 | 文 |
| 関係抽出 | 20,692 | 文 |
| 同一性判定 | 約 3,000 万 | 文 |
| 分類 | 2,000 | アブストラクト |

がある。例えば、キュレーターが修正した通りに抽出できるようになっているのか、修正箇所だけでなく全体として精度が上がっているのか、などといった検証を行う必要がある。

5. 実験

本章では、キュレーターの作業効率化について考察するために、開発した文献キュレーションシステムの文献検索やアノテーションについて実験する。また実験を通す中で、UIについても評価を行う。

5.1 実験データ

実験に使用する学習データと文献データは、PubMed[15]から取得した論文データを使用する。PubMedとは、アメリカ国立医学図書館が提供する医療文献検索サービスである。MEDLINE[16]やPMC[17]といった文献データベースを対象として、オープンに検索できる。

表1に学習データの内訳を示す。固有表現抽出や関係抽出など、処理内容によって学習データも異なる。これらの学習データから学習モデルを構築し、PMCから取得した

表 2 実験用 CIViC データ

| 実験データ No. | Protein | Mutation | Disease | Drug | Type | PubMedID |
|-----------|---------|----------|------------------------------|-----------|------------|----------|
| 1 | BTK | T316A | Chronic Lymphocytic Leukemia | Ibrutinib | Resistance | 27626698 |
| 2 | HOXB13 | G84E | Prostate Cancer | - | Positive | 27626483 |

表 3 検索条件と結果

| 検索条件 No. | 実験データ No. | 固有表現タイプ | キーワード 1 | キーワード 2 | 関係 | 結果 |
|----------|-----------|------------------------|------------------------------|-----------------|----------|--------|
| 1 | 1 | Mutation - Drug | T316A | Ibrutinib | Negative | 該当論文なし |
| 2 | 1 | Mutation - Drug | T316A | - | Negative | 該当論文なし |
| 3 | 1 | Mutation - Drug | - | Ibrutinib | Negative | 該当論文あり |
| 4 | 1 | Gene.Protein - Drug | BTK | Ibrutinib | Negative | 該当論文なし |
| 5 | 1 | Gene.Protein - Drug | BTK | - | Negative | 該当論文なし |
| 6 | 1 | Gene.Protein - Drug | - | Ibrutinib | Negative | 該当論文なし |
| 7 | 1 | Disease - Drug | Chronic Lymphocytic Leukemia | Ibrutinib | Negative | 該当論文なし |
| 8 | 1 | Disease - Drug | Chronic Lymphocytic Leukemia | - | Negative | 該当論文なし |
| 9 | 1 | Disease - Drug | - | Ibrutinib | Negative | 該当論文なし |
| 10 | 2 | Disease - Mutation | Prostate Cancer | G84E | Positive | 該当論文なし |
| 11 | 2 | Disease - Mutation | Prostate Cancer | - | Positive | 該当論文なし |
| 12 | 2 | Disease - Mutation | - | G84E | Positive | 該当論文あり |
| 13 | 2 | Gene.Protein - Disease | HOXB13 | Prostate Cancer | Positive | 該当論文なし |
| 14 | 2 | Gene.Protein - Disease | HOXB13 | - | Positive | 該当論文なし |
| 15 | 2 | Gene.Protein - Disease | - | Prostate Cancer | Positive | 該当論文なし |

約 13 万件の論文を文献データとして医療情報を抽出した。

5.2 文献検索の実験

文献検索画面において、固有表現・関係について検索リクエストし、適切な文献が結果として表示されるか実験する。適切かどうかの判断は、キュレーション済みの医療情報が登録されたデータベースである CIViC[18] を利用する。表 2 は実験に使用した CIViC のデータである。例えば一行目（実験データ No.1）は、治療薬 Ibrutinib が疾患 Chronic Lymphocytic Leukemia に有効でないことを示す論文として、PubMed ID:27626698 があることを意味している。なお、PubMed ID とは PubMed で検索できる文献に振られた ID である。

表 3 は、実験データに対する検索条件と検索結果である。CIViC の実験データに合わせて検索条件を指定しており、各実験データに対し、条件を変えながら複数回検索している。結果が「該当論文あり」であれば、実験データに含まれる PubMed ID が検索結果に現れたことを示し、適切な結果が得られたことになる。

結果としては、実験データ No.1 については、検索条件 No.3 において該当論文を得られた。しかし、それ以外の検索条件については該当論文を得られなかった。実験データ No.2 についても似たような結果となり、検索条件 No.12 において該当論文を得られたが、それ以外の検索条件では得られなかった。今回の実験データにおいては、検索精度は決して高くはなく改善が必要と言えよう。

5.3 アノテーションの実験

次に文献キュレーション画面において、適切にアノテーションが付与されているか実験する。実験した論文の PubMed ID を以下に示す。

- 30076320
- 30083449
- 29189985

各論文の文献キュレーション画面を表示し、医療文献キュレーションの有識者によってアノテーションが適切かどうか判断する。

PubMed ID:30076320 の結果を図 8 に示す。赤丸は誤っていると判断されたアノテーションである。逆に赤丸がなければ適切と判断されたアノテーションである。例えば一行目においては特に赤丸はないため、「geomoitabine」は「薬剤名」、「breast cancer」は「疾患名」、それが有効であることが適切であると示している。一方、TQ という固有表現は赤丸が付いているが、これは「Drug」ではなく、正しくは「Mutation」となり誤りである。図中において、誤りは 3 箇所あるが、正解は関係まで含めれば 31 箇所あるため、概ね正しく抽出できている。アノテーションがない場合は 33 箇所全てにアノテーション作業が発生するが、文献キュレーションシステムの場合は 3 箇所のみ修正となるため、作業効率は向上すると言える。

5.4 UI の評価

本実験を行ったキュレーションの有識者によって、実験

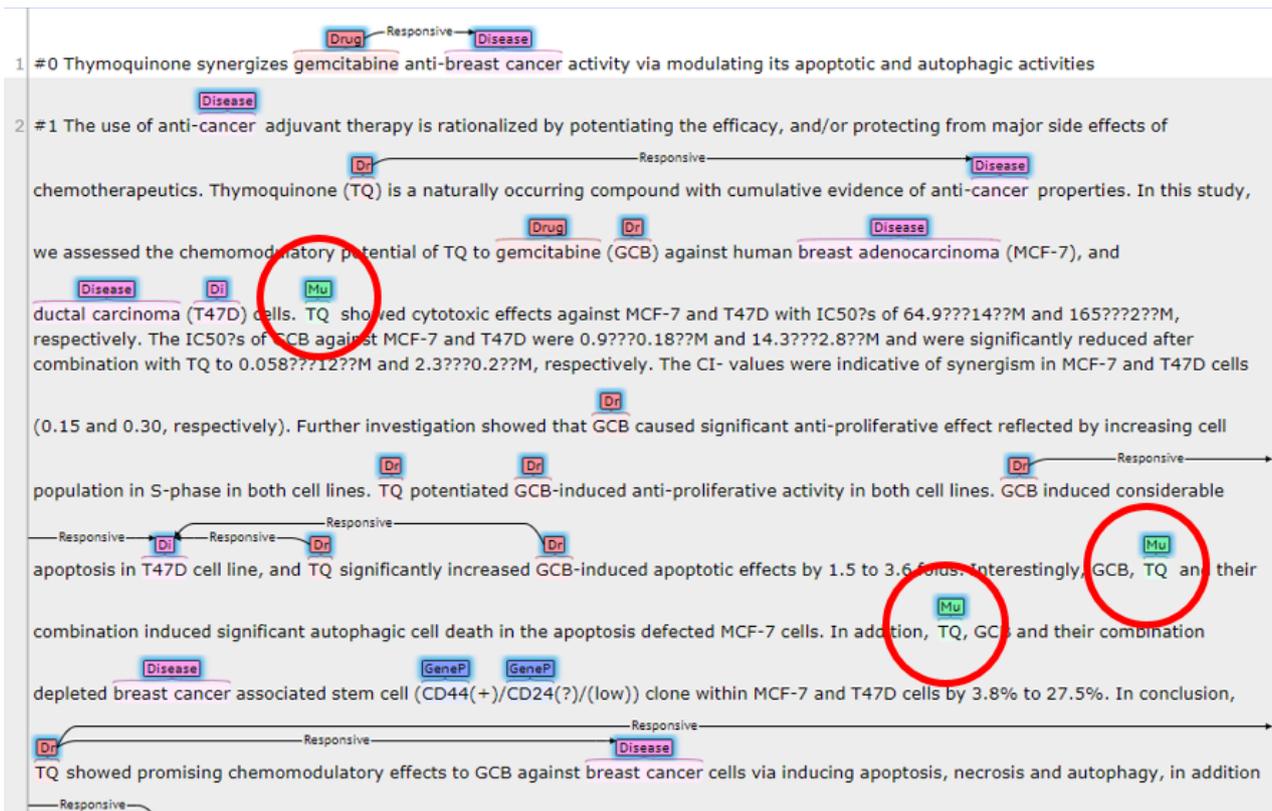


図 8 アノテーション実験結果

を通す中で UI についても評価を行った。それにより、抽出した課題と改善案を以下に示す。

複数の検索条件のサポート

現在の文献検索画面は、一つの固有表現・関係のセットしか検索できない。一般的な検索でもあるように、複数の条件で AND 検索や OR 検索を行うことで、より検索者が要求する文献に絞り込める。そのため、複数の固有表現・関係セットで検索を行えるように、検索フォームを改善する。

オントロジーの導入

現在はキーワードマッチで検索しており、キーワードの意味を考慮していない。意味を考慮した検索を行う方法として、オントロジーの導入が挙げられる。オントロジーとは、ある特定分野における体系化された概念や知識を指す。例えば、Cancer の下位概念として Lung Cancer があり、さらに Lung Cancer の下位概念として NSCLC や Small-cell lung cancer がある。このようなオントロジーを文献検索画面にて図示し、概念を選択することで検索リクエストを行えるようにする。それにより、例えば上位概念を選択すると下位概念全てが検索対象とするような、より意味的に正確な検索結果を期待できる。

6. まとめと今後の課題

本研究では、医療文献のキュレーション作業効率化のため

に、アノテーションツールと自然言語処理技術を連携させた文献キュレーションシステムを開発した。いくつかの文献で実験したところ、自然言語処理によって大部分のアノテーションは正しく自動抽出しており、キュレーション作業の効率化が図れることを確認した。文献検索においては、文献が持つ分類情報や関係情報を提示することで、文献の選択がしやすくなった。しかし、文献検索の実験では、検索条件にマッチするような適切な文献を提示できなかったケースが多い。またフィードバックについては、今回は実験を行っていないが、フィードバックのタイミング、頻度、再学習後の精度検証など課題も多い。

従って今後は、オントロジーの導入による検索精度の改善や、フィードバックの機能強化を行う。そして、全体を通してキュレーション作業がどれほど効率化されるか定量的に評価を行う予定である。

謝辞 本研究の一部は AMED の課題番号 JP18kk0205013 の支援を受けた。

参考文献

- [1] Neves, M. and Leser, U.: A survey on annotation tools for the biomedical literature, *Briefings in bioinformatics*, Vol. 15, No. 2, pp. 327–340 (2012).
- [2] Kilicoglu, H., Demner-Fushman, D., Rindflesch, T. C., Wilczynski, N. L. and Haynes, R. B.: Towards automatic recognition of scientifically rigorous clinical research evidence, *Journal of the American Medical Informatics Association*, Vol. 16, No. 1, pp. 25–31 (2009).

- [3] Fiszman, M., Demner-Fushman, D., Kilicoglu, H. and Rindfleisch, T. C.: Automatic summarization of MEDLINE citations for evidence-based medical treatment: A topic-oriented evaluation, *Journal of biomedical informatics*, Vol. 42, No. 5, pp. 801–813 (2009).
- [4] Wallace, B. C., Trikalinos, T. A., Lau, J., Brodley, C. and Schmid, C. H.: Semi-automated screening of biomedical citations for systematic reviews, *BMC bioinformatics*, Vol. 11, No. 1, p. 55 (2010).
- [5] Cohen, A. M., Smalheiser, N. R., McDonagh, M. S., Yu, C., Adams, C. E., Davis, J. M. and Yu, P. S.: Automated confidence ranked classification of randomized controlled trial articles: an aid to evidence-based medicine, *Journal of the American Medical Informatics Association*, Vol. 22, No. 3, pp. 707–717 (2015).
- [6] Sarker, A., Mollá, D. and Paris, C.: Automatic evidence quality prediction to support evidence-based decision making, *Artificial intelligence in medicine*, Vol. 64, No. 2, pp. 89–103 (2015).
- [7] Institute for Genome Sciences: Disease Ontology, <http://disease-ontology.org/>. (accessed 2019-02-22).
- [8] University of Alberta and The Metabolomics Innovation Centre: DrugBank, <https://www.drugbank.ca/>. (accessed 2019-02-22).
- [9] HGNC: HUGO Gene Nomenclature Committee, <https://www.genenames.org/>. (accessed 2019-02-22).
- [10] Stenetorp, P., Pyysalo, S., Topić, G., Ohta, T., Ananiadou, S. and Tsujii, J.: BRAT: a web-based tool for NLP-assisted text annotation, *Proceedings of the Demonstrations at the 13th Conference of the European Chapter of the Association for Computational Linguistics*, Association for Computational Linguistics, pp. 102–107 (2012).
- [11] Garfield, E.: The history and meaning of the journal impact factor, *Jama*, Vol. 295, No. 1, pp. 90–93 (2006).
- [12] Priem, J., Taraborelli, D., Groth, P. and Neylon, C.: Altmetrics: A manifesto (2010).
- [13] Hirsch, J. E.: An index to quantify an individual’s scientific research output, *Proceedings of the National Academy of Sciences*, Vol. 102, No. 46, pp. 16569–16572 (2005).
- [14] CEBM: Oxford Centre for Evidence Based Medicine, <https://www.cebm.net/>. (accessed 2019-02-22).
- [15] National Center for Biotechnology Information: PubMed, <https://www.ncbi.nlm.nih.gov/pubmed/>. (accessed 2019-02-22).
- [16] National Center for Biotechnology Information: MEDLINE, <https://www.nlm.nih.gov/bsd/medline.html>. (accessed 2019-02-22).
- [17] National Center for Biotechnology Information: PMC, <https://www.ncbi.nlm.nih.gov/pmc/>. (accessed 2019-02-22).
- [18] The McDonnell Genome Institute: CIViC, <https://civicdb.org/home>. (accessed 2019-02-22).