

科学研究費助成事業（科学研究費補助金）研究成果報告書

平成 25 年 6 月 7 日現在

機関番号：12601

研究種目：挑戦的萌芽研究

研究期間：2010 ～ 2012

課題番号：22650025

研究課題名（和文） タミー診療録の構築および自動構造化に関する研究

研究課題名（英文） Automatic Structuring for Dummy Health Records

研究代表者

荒牧 英治 (ARAMAKI EIJI)

東京大学・知の構造化センター・講師

研究者番号：70401073

研究成果の概要（和文）：

平成 13 年度に政府が発表した「保健医療分野の情報化にむけてのグランドデザイン」にて、電子カルテシステムを始め医療 IT 技術の普及が課題の一つとして掲げられました。以降、急速に医療の IT 化が進み、その結果、かつてない大量の臨床データが電子化された状態でストックされつつあります。このデータを有効に利用することができれば患者の生活習慣と疾患の相関（例えば、喫煙と癌）や、薬品と副作用の相関（タミフルと精神障害）について過去類を見ない大規模な調査を迅速に行うことが可能となり、臨床研究が加速的に進展するとして高い期待が寄せられています。しかし、単にデータを電子化しただけで、上記のような革新が実現できるわけではありません。電子カルテにおいても自然言語で入力される箇所が相当な割合で存在します。よって、電子カルテデータを臨床知識としてフルに活用するためには、自然言語処理を活用することが必須となります。以上の背景のもと、本プロジェクトでは模擬診療録を用いた応用研究を進めます。

研究成果の概要（英文）：

The use of Electronic Health Records (EHR) in hospitals is increasing rapidly everywhere. They contain much clinical information about a patient's health, including the frequency of drug usage, related side-effects, and so on, which facilitates unprecedented large-scale research. Nevertheless, extracting clinical information from the reports is not easy because they are written in natural language. This study specifically examines clinical information using dummy records that are also build by this project.

交付決定額

(金額単位：円)

	直接経費	間接経費	合計
2010 年度	1,400,000	0	1,400,000
2011 年度	1,000,000	300,000	1,300,000
2012 年度	500,000	150,000	650,000
年度			
年度			
総計	2,900,000	450,000	3,350,000

研究分野：知識処理

科研費の分科・細目：知識処理・自然言語処理

キーワード：自然言語処理, 医療情報学

1. 研究開始当初の背景

平成 13 年度に政府が発表した「保健医療分野の情報化にむけてのグランドデザイン」にて、電子カルテシステムを始め医療 IT 技術の普及が課題の一つとして掲げられた。以降、急速に医療の IT 化が進み、その結果、かつてない大量の臨床データが電子化された状態でストックされつつある。このデータを有効に利用することができれば患者の生活習慣と疾患の相関（例えば、喫煙と癌）や、薬品とその副作用の相関（タミフルと精神障害）について過去類を見ない大規模な調査が可能となり、臨床研究が加速的に進展するとして高い期待が寄せられている。

しかし、単にデータを電子化しただけで、上記のような革新が実現できるわけではない。この理由の一つは、電子カルテにおいても自然言語で入力される箇所が相当な割合で存在するためである。よって、電子カルテデータを臨床知識としてフルに活用するためには、自然言語情報を活用することが必須となる。

2. 研究の目的

本研究では言語処理技術を用い、カルテのテキスト情報を構造化された情報へ変換する技術を研究／開発し、その技術で臨床現場をサポートするシステムを構築／実証実験を行う。

3. 研究の方法

カルテ文章は個人情報の塊ともいえる文章であり、研究利用が困難である。したがって、まず、ダミーのカルテデータを構築する (phase-1)。(図 1 参照)

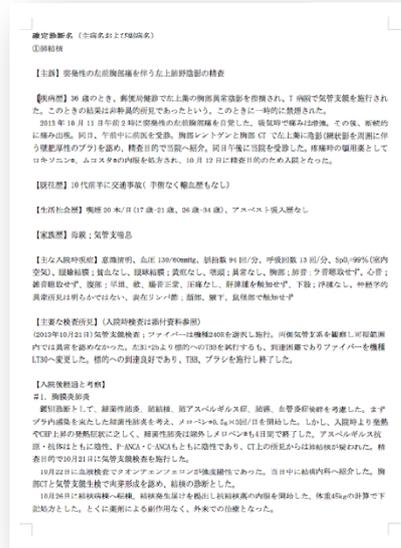


図 1：ダミーカルテの例 (平均的な 1 枚のカルテ)

工場に勤めている<a>64歳の<x>男性</x>。<t>2025
年8月2日(未病5日前)頃から<t><c>腹痛</c>が生じると
ともに、<c>食欲不振</c>、<c>嘔気</c>、<c>嘔吐出現</
c>した。体幹は温かいが、末梢は<c>湿潤冷汗</c>で<c>
ショック状態</c>。明らかな<c modality="negation">
運動麻痺</c>はみられず。<t>翌日</t>、<c>意識障害出
現</c>し、<c>腎機能障害</c>の増悪を認めて徐々に<c>
尿量低下</c>し、<t>8月9日18時10分</t>に<c>心肺停止
</c>。<t>8月9日21時44分</t><c>死亡確認</c>。

図2：ダミーカルテのアノテーション例

次にデータのアノテーションを行い
(phase-2)、これを自動化する言語処理技術
を研究開発する。研究対象となる技術は、固
有表現特定/表記ゆれ吸収/モダリティ解
析といった基礎研究に分解可能で、学術的意
義と社会的重要性の両方に貢献可能なよう
配慮する。

最後に、開発した技術を用いたアプリケー
ションを実装し、実験を行う。

実験は、固有表現の自動特定、症状の特定
などを行う。

実験の結果、表1のような精度が得られた。

tag	Precision	Recall	$F_{\beta=1}$
<a>: age	88.89	75.00	81.36
<t>: time	83.08	76.60	79.70
<h>: hospital	96.15	65.79	78.12
<l>: location	0.00	0.00	0.00
<x>: sex	100.00	50.00	66.67
<c>	87.37	71.86	78.86
<c> positive	62.58	62.08	62.33
<c> family	100	59.09	74.29
<c> negation	78.47	66.4	71.93
<c> suspicion	60	20	30

表1：実験結果

age は年齢, time は日付, hospital は医院
名, location は地域名, sex は性別を示す <c>
は症状名を示す。

図3：匿名化コピー

また、このシステムを実装したハード(匿名
化コピー)を開発し、発表を行った。この匿
名化コピーは多くのメディアに取り上げら
れ商品化を複数社と検討している。

4. 研究成果

本研究の結果、50枚のダミーカルテ文章が完
成した。このダミーカルテを研究目的にて配
布した。また、ダミーカルテをもとに国際ワ
ークショップ NTCIR10 MedNLP を主催した。
これは本邦ではじめてのカルテデータを用
いたシェアードタスクであり、国外を含め16
カ国の研究チームが参加した。さらに、2015
年度の開催も決定し、今後の発展が期待され
る。

5. 主な発表論文等

[雑誌論文] (計0件)

[学会発表] (計7件)

(1) Yasuhide Miura, Tomoko Ohkuma,
Hiroshi Masuichi, Emiko Yamada Shinohara,
Eiji Aramaki, Kazuhiko Ohe: UT-FX at
NTCIR-10 MedNLP: Incorporating Medical
Knowledge to Enhance Medical Information
Extraction, In Proceedings of NTCIR-10,
2013. (2013年6月18日, 東京)

(2) Mizuki Morita, Kano Yoshinobu,
Ohkuma Tomoko, Miyabe Mai, Eiji Aramaki:
Overview of the NTCIR-10 MedNLP task, In
Proceedings of NTCIR-10, 2013. (2013年6
月18日, 東京)

(3) Hiroto Imachi, Mizuki Morita, Eiji Aramaki: NTCIR10 MedNLP Baseline System, In Proceedings of NTCIR-10, 2013. (2013年6月18日, 東京)

(4) 森田瑞樹, 狩野芳伸, 大熊智子, 宮部真衣, 荒牧英治: NTCIR-10 “MedNLP” Pilot Task: 医療分野の言語処理研究の環境整備に向けて, 言語処理学会 第19回年次大会, 2013. (2013年3月13-15日)

(5) 森田瑞樹, 荒牧英治: 医療分野における言語処理研究の環境整備に向けての提案,

テキスト
ノテーショ
ンワークシ
ョップ
2012. (2012
年8月6日,
東京).
(奨励賞;
16%=3件/18
件)



(6) 荒牧英治, 増川佐知子, 森田瑞樹, 保田祥: 日本人のオンライン・コミュニケーション上での平均使用語彙数は8,000語である, 情報処理学会 第208回自然言語処理研究会 (SIG-NL), 2012. (2012年9月3日, 仙台)

(7) 森田瑞樹, 篠原(山田)恵美子, 宮部真衣, 荒牧英治: 受診理由のコーディング支援ツールの開発, 第32回医療情報学連合大会, 2012. (2012年11月16日, 新潟)

[図書] (計0件)

[産業財産権]

○出願状況 (計0件)

名称:
発明者:
権利者:
種類:
番号:
出願年月日:
国内外の別:

○取得状況 (計0件)

名称:
発明者:
権利者:
種類:
番号:
取得年月日:

国内外の別:

[その他]

<http://mednlp.jp>

6. 研究組織

(1) 研究代表者

荒牧 英治 (ARAMAKI EIJI)

東京大学・知の構造化センター・講師

研究者番号: 70401073