

医療テキストデータに対するトピックモデルの応用

大谷 誠¹⁾、松田 晋哉^{1,2)}

1) 産業医科大学 産業保健データサイエンスセンター

2) 産業医科大学 医学部 公衆衛生学教室

要旨：本研究では医療データである DPC データを用いテキストマイニングを行った。医療データは数値データやカテゴリー尺度、ICD10 といったコードなど分析し医療行為にフィードバックが行えるよう精緻化されてきている。DPC コード毎のトピックの出現度をみると、手術・処置有無に関してリハビリは前回手術を行った DPC コード、症状不良は前回入院時に生検法をした化学療法と放射線治療の両方を行っている DPC コードに関連していた。また、化学療法と放射線治療の両方が行われている DPC コードは 3 から 4 つのトピックを含んでいた。本研究によりテキスト分析により医療行為内容の体系的把握が可能になる可能性が示された。医療行為の均一化や医療職の負担軽減を実現するために、より精度の高いテキスト分析の手法を開発していく必要がある。

キーワード： DPC データ、テキストマイニング、トピックモデル

1. はじめに

平成 30 年 8 月戦略的イノベーション創造プログラム第 2 期において AI ホスピタルによる高度診断・治療システムが掲げられている。AI ホスピタルによるシステムは「高度で先進的な医療サービスを提供するとともに、医療機関における効率化を図り、医師や看護師などの医療従事者の抜本的な負担の軽減を実現する。」ものである。システム構築に向け 2020 年度までの到達目標の 1 つに「セキュリティの高い医療情報データベースの構築・医療有用情報抽出技術の開発」が掲げられている¹⁾。

医療情報には血液検査やバイタルサインなどの数値データ、CT や MRI などの画像データ、SOAP などで電子カルテに記入されるテキストデータなどがある。画像データに関しては G. Litjens の論文によると 2016 年以降に Deep Learning を用いた画像解析に関する論文等が発表されていると報告されている²⁾。

情報抽出の技術としてデータマイニングがある。データマイニングとは「データの中に埋め込まれている知識の発掘」である³⁾。データマイニングの分析対象がテキストデータとなったものがテキストマイニングである。医療分野でのテキストマイニングは 1960 年代から行われている⁴⁾。日本では退院サマリーからの疾患名の特定⁵⁾や看護記録からの有用情報の抽出⁶⁾に関する研究が行われている。医療分野でテキストデータが入力されているものとして DPC (Diagnosis Procedure Combination) データがある。

DPC データは患者の基本情報や入院先に関する情報、疾患の重症度など診療録情報が記載される様式 1、実施・処方される手術や処置、薬剤の情報など診療報酬請求情報が記載される入院 EF 統

合ファイルなどで構成されている⁷⁾。様式1には再入院および再転棟に関する自由記載の入力欄がある。本研究では再入院に着目しそのうち計画外の再入院のテキストデータから再入院する患者の傾向をテキストマイニングで抽出することを目的に分析を行う。

2. 方法

2.1 分析データおよび分析対象

一般社団法人診断群分類研究支援機構が収集した平成26年度及び平成27年度のDPCデータを用いた。分析対象としてDPC6桁コード(病名)が「肺の悪性腫瘍(040040)」の患者を選択した。DPCデータの再入院には計画的再入院、予期された再入院、予期せぬ再入院の3種類が用意されている。本研究では「予期された再入院および予期せぬ再入院」の患者を抽出し、自由記述欄に記載がある患者(42,605名)を対象とした。

2.2 分析方法

本研究ではテキストマイニングの手法にトピックモデルを用いた。トピックモデルとは文書に出現する単語とその出現回数の情報からそれぞれの文書に潜在的に存在するトピックを推定する手法である⁸⁾。本研究ではトピックモデルの手法の1つである潜在的ディレクレ配分法(Latent Dirichlet Allocation: LDA)⁹⁾を用いた。トピックモデルによるテキストマイニングの概要を図1にLDAのトピック生成過程のアルゴリズムを以下に示す。本研究ではトピック数を $K=7$ とした。

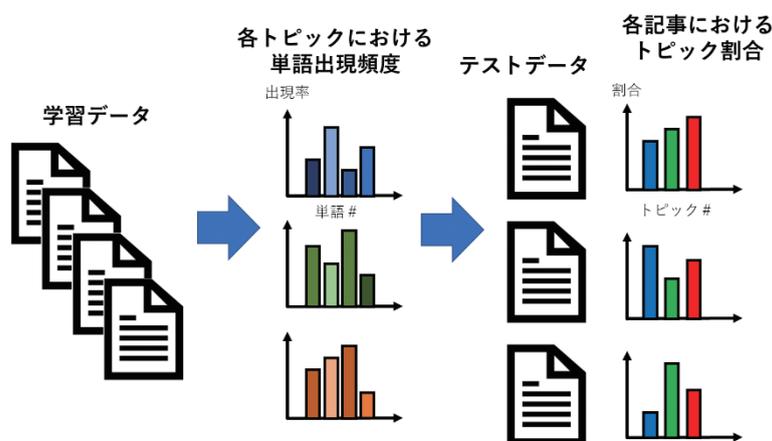


図1 トピックモデルの概要

2.3 学習データ

トピックモデルに学習させるデータの作成は以下の手順で行った(図2参照)。

- ① 入院患者毎に前回入院時のDPCコード(14桁)と再入院理由の自由記載を抽出
- ② DPCコード毎に自由記載を結合(前処理として句読点や記号の除去、半角全角や大文字小文字の統一を行った)
- ③ 結合後の文章に対し形態素解析を行い名詞(一般、固有名詞、形容動詞語幹)と形容詞(自立)の単語を抽出(形態素解析にはPythonのパッケージであるJanomeを用いた)

アルゴリズム 1 生成アルゴリズム

```

for トピック  $k = 1, \dots, K$  do
  単語分布を生成  $\phi_k \sim \text{Dirichlet}(\beta)$ 
end for
for 文書  $d = 1, \dots, D$  do
  トピック分布を生成  $\theta_d \sim \text{Dirichlet}(\alpha)$ 
  for 単語  $n = 1, \dots, N_d$  do
    トピックを生成  $z_{d,n} \sim \text{Categorical}(\theta_d)$ 
    単語を生成  $w_{d,n} \sim \text{Categorical}(\phi_{z_{d,n}})$ 
  end for
end for
  
```

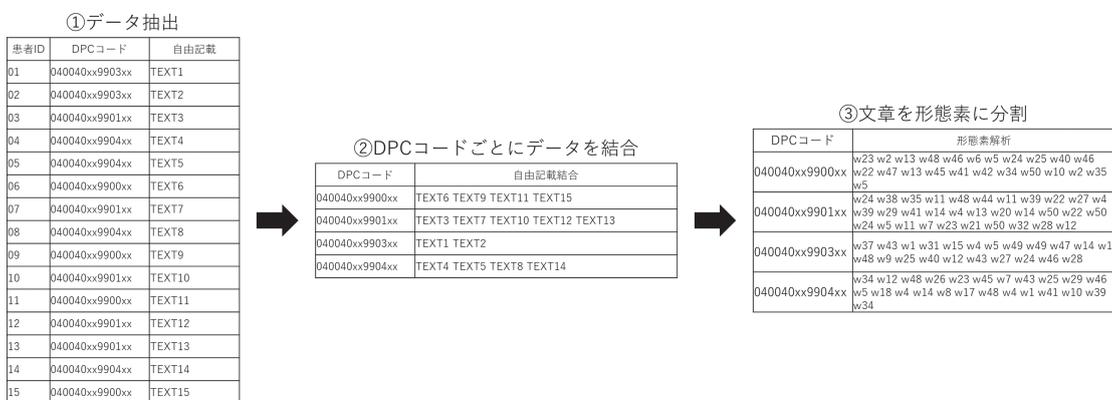


図 2 学習データ作成フロー図

3. 結果

トピックモデルによる各トピックと単語の関係を図3および表1に示す。図3は各トピックで出現する単語の頻度を文字の大小で表現したものである。表1は各トピックに対して出現する単語からトピック名をラベル付けしたものである。DPCコード毎のトピックの出現度を表2に示す。Topic#2のりハビリに関しては前回入院時に手術を行ったコードのみ表れた。Topic#6の症状不良に関しては前回入院時に生検法をされており化学療法と放射線治療の両方を行っているコードのみに表れた。また化学療法と放射線治療の両方が行われているコードは3から4つのトピックが含まれていた(040040xx97x4xx, 040040xx9904xx, 040040xx9914xx)。Topic#5の治療に関しては全コード中の約8割、Topic#0の検査に関しては約4割で現れた。

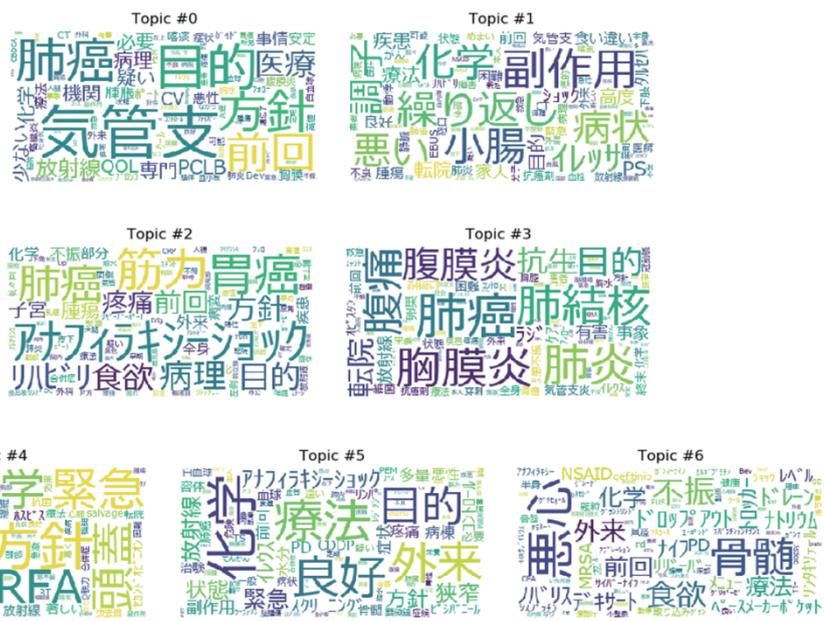


図3 トピック内の単語出現頻度

表1 トピック名

Topic#	トピック名	係数*単語
0	検査	0.009* 気管支 +0.007* 目的 +0.006* 肺癌 +0.006* 方針 +0.006* 前回 +0.005* 医療 +0.005* 放射線 +0.005*PCLB+0.004* 必要 +0.004* 検査
1	治療不良	0.007* 副作用 +0.006* 繰り返し +0.006* 小腸 +0.005* 化学 +0.005* 病状 +0.005* 悪い +0.005* 調子 +0.004* イレッサ +0.004* 疾患 +0.004* 治療不良
2	リハビリ	0.006* アナフィラキシーショック +0.005* 筋力 +0.005* 胃癌 +0.004* 肺癌 +0.004* 病理 +0.004* リハビリ +0.004* 目的 +0.004* 食欲 +0.004* 方針
3	他疾患	0.007* 肺癌 +0.006* 肺炎 +0.006* 胸膜炎 +0.005* 肺結核 +0.005* 腹痛 +0.005* 腹膜炎 +0.005* 目的 +0.005* 抗生 +0.005* 転院 +0.004* 他疾患
4	緊急	0.004* 方針 +0.004* 目的 +0.004* 救急 +0.003* 緊急 +0.003* 化学 +0.003* 頭蓋 +0.003*RFA+0.003* ロキソプロフェン +0.003* 内圧
5	治療	0.010* 化学 +0.008* 療法 +0.008* 良好 +0.007* 外来 +0.007* 目的 +0.005* アナフィラキシーショック +0.005* 緊急 +0.005* 狭窄 +0.005* 放射線 +0.005* 治療
6	症状不良	0.009* 悪心 +0.006* 骨髄 +0.003* 不振 +0.003* 食欲 +0.002* 外来 +0.002* 前回 +0.002* 化学 +0.002* 療法 +0.001* ドレーン +0.001* 症状不良

表2 DPCコードとトピックの関係

DPC14	Topic #0	Topic #1	Topic #2	Topic #3	Topic #4	Topic #5	Topic #6	手術	処置 1	処置 2
040040xx97x0xx	0.37		0.12	0.02		0.49		○	×	×
040040xx97x10x				0.86		0.14		○	×	○
040040xx97x2xx	0.90			0.10				○	×	○
040040xx97x3xx		0.56				0.44		○	×	○
040040xx97x4xx	0.05	0.05			0.17	0.72		○	×	○
040040xx97x5xx	1.00							○	×	○
040040xx97x6xx			0.57			0.43		○	×	○
040040xx97x7xx						1.00		○	×	○
040040xx97x80x						1.00		○	×	○
040040xx97x81x						1.00		○	×	○
040040xx9900xx	0.18			0.39		0.43		×	×	×
040040xx9901xx				0.99				×	×	○
040040xx9902xx	0.26				0.41	0.34		×	×	○
040040xx9903xx						1.00		×	×	○
040040xx9904xx	0.03	0.01		0.05		0.91		×	×	○
040040xx9905xx		0.73				0.27		×	×	○
040040xx9906xx						1.00		×	×	○
040040xx9907xx		0.29				0.70		×	×	○
040040xx9908xx	0.41					0.59		×	×	○
040040xx99100x	0.89					0.11		×	○	×
040040xx99101x	0.78					0.22		×	○	×
040040xx9911xx	1.00							×	○	○
040040xx9912xx				0.68		0.32		×	○	○
040040xx9913xx						0.68	0.31	×	○	○
040040xx9914xx	0.22	0.43				0.36		×	○	○
040040xx9915xx						1.00		×	○	○
040040xx9916xx				0.51		0.49		×	○	○
040040xx9917xx		1.00						×	○	○

4. 考察

テキストマイニングの対象として再入院理由を選択した。再入院率は医療の質の評価指標¹⁰⁾として用いられており、再入院率を下げるのが期待されている。本研究では計画的再入院以外の再入院を分析対象とした。現在、治療と仕事の両立支援¹¹⁾が国から各企業に求められており就労者にとっても計画的に治療が進むことで就労への障害が減りより積極的な治療が望めると思われる。本研究では肺

の悪性腫瘍を対象とした。2018年度より診療報酬としてがん患者対象に療養・就労両立支援指導料が始まったこと、がんの中の罹患者数が3番目に多く死亡数が一番多い¹²⁾ことから対象の疾患とした。再入院にあたり前回入院における医療行為により理由のパターンがあると仮説を立てDPC14桁別に分析を行ったがパターンは見つけることができなかった。その要因として2つ考えられ①用語のばらつきと②各施設の入力内容の粒度の違いが考えられる。①に関しては入力される用語は日本語入力や英語入力、略語などバリエーションがあった。例えばリハビリテーションは「リハビリテーション」、「リハビリ」、「リハ」である。この対策としては用語の対応表を作成し正式名称に変換後、分析する必要があると考えられる。②に関しては丁寧に入力している施設もある一方、2文字しか入力していない施設があった。データを作成したテキストは自由記載欄ではあるが、知識の共有や分析への応用を考えると最低でも「文」として入力を求める必要があると思われる。本研究ではトピックに割り当てられた単語から研究者がトピック名の決定を行った。トピック名決定には専門的な知識が必要でありデータによって異なる専門家が必要になると思われる。

本研究では医療データであるDPCデータを用いテキストマイニングを行った。医療データは数値データやカテゴリー尺度、ICD10といったコードなど分析し医療行為にフィードバックが行えるよう精緻化されてきている。しかし電子カルテ内にはテキストデータが多く存在していることも事実でありその中に共有すべきノウハウが存在している。そのためテキストからノウハウを抽出し共有することで医療行為の均一化や医療職の負担軽減に寄与すると考えられることから、より精度の高いテキスト分析の手法を開発していく必要があると思われる。

参考文献

- 1) 内閣府, 戦略的イノベーション創造プログラム (SIP) の概要 研究開発計画 第2期 12 課題, <https://www8.cao.go.jp/cstp/gaiyo/sip/kenyugaiyo2.pdf>, 2019年12月10日閲覧
- 2) Geert Litjens, et al., A survey on deep learning in medical image analysis, *Medical Image Analysis*, Vol. 42, pp. 60-88, 2017
- 3) 人工知能学会編, 人工知能学辞典, 協立出版株式会社, 2005
- 4) Hercules Dalianis, *Clinical Text Mining Secondary Use of Electronic Patient Records*, Springer, 2018
- 5) 小野大樹ら, テキストマイニングによる退院サマリー自動分類の試み, *医療情報学*, 24(1), pp 35-44, 2004
- 6) 村松洋, 看護記録のテキストマイニング, *情報処理学会論文誌データベース*, Vol. 3, No. 3, pp. 112-122, 2010
- 7) 平成30年度「DPC導入の影響評価に係る調査」実施説明資料, https://www.prrism.com/dpc/setumei_20180406.pdf, 2019年12月10日閲覧
- 8) 田村一樹ら, 評点付きレビュー文書を対象としたトピックモデルの構築に関する検討, *情報処理学会論文誌*, Vol. 56, No. 3, pp. 1013-1027, 2015
- 9) David M. Blei, Latent Dirichlet Allocation, *Journal of Machine Learning Research*, 3, pp. 993-1022, 2003
- 10) 厚生労働省, 平成30年度医療の質の評価・公表等推進事業公募要領 共通指標セット (参考資料1), https://www.mhlw.go.jp/file/06-Seisakujouhou-10800000-Iseikyoku/0000166398_2.pdf,

2019年12月10日閲覧

- 11) 厚生労働省, 事業場における治療と仕事の両立支援のためのガイドライン,
<https://www.mhlw.go.jp/content/11200000/000490701.pdf>, 2019年12月10日閲覧
- 12) 国立がん研究センター, がん統計 最新がん統計, https://ganjoho.jp/reg_stat/statistics/stat/summary.html, 2019年12月10日閲覧