

LLMを用いた呼吸器内科SOAPデータの分類、特徴抽出及び構造化

張 孜恒¹、韓 豊¹、唐 明輝²、平田 健司²、北井秀典³、小笠原 克彦⁴、工藤 興亮²

- 北海道大学大学院医学院 画像診断学教室
- 北海道大学大学院医学研究院 画像診断学教室
- 北海道大学大学院医学研究院 呼吸器内科学教室
- 北海道大学大学院保健科学研究所 健康科学分野

Introduction & Objective

電子カルテ (EMR) の普及により、SOAP (Subjective, Objective, Assessment, Plan) 形式の診療記録が標準的に使用されている^[1]。しかし、SOAPは記述の自由度が高く、情報の標準化や自動処理が難しいという課題がある^[2]。

例えば、同じ診断・症状でも記載方法が異なる場合が多く、「主訴 (CC)」や「現在の病気の履歴 (HPI)」の境界が曖昧になったり、「診断 (Assessment)」の記述が治療計画 (Plan) に含まれることもある。また、医師ごとの表記の違いにより、自動処理やデータ解析の妨げとなる^[3]。そのため、SOAPフォーマット内の情報を、**統一したルールに基づいて詳細なカテゴリ (第一層) とその細分類 (第二層) に分け、標準化・構造化することが重要である**

本研究では、日本語の呼吸器内科SOAPデータを対象に、大規模言語モデル (LLM) を用いた自動アノテーション手法を開発・評価し、診療記録の情報抽出と構造化の精度を向上させることを目的とする。

Methods

1. データセット

200例の呼吸器内科SOAPデータ (医師による確認済み、ゴールドスタンダードデータとして使用)

- トレーニング用データ160例
- 評価用データ40例

S(症状)	昨日よりずいぶん体が楽である。本日、洗濯後排便あり、とても楽になった。以前は毎朝8時30分に排便があったが、ここ1週間は排便なかった。少しからだを動かそうとして室内を歩いたら、看護師さんの止められて、今日は安静にしようと思う。
O(所見)	BT37.8-36.9度 P110 BP130/81 SpO2 98%(room air) BW40.50
症例数:	L/D
11万2420	01 白血球数 5.6 緊
	02 赤血球数 3.13 L 緊
	03 ヘモグロビン 8.5 L 緊
病名:	...
1. 肺炎	30 カルシウム 8.1 L 緊 ZT
2. 肺腫瘍	
3. 気管支喘息	CXp (ポーター臥位)
4. 肺線維症	
5. 肺炎	
A(判断)	コントラストの影響もあるが8/3と比較し両肺野の透過性は改善 昨日、入院時と比較し呼吸状態、全身状態の悪化なし 貧血は赤血球輸血2単位によりHb8.5まで改善 低Na進行注意 廃用症候群注意→明日以降、付き添いで機内フリー等を御検討いただく
P(計画)	アザクテム、ダランシによる加療継続 明日、CXp、L/Dにて再評価

Figure 1: Example of SOAP data

2. アノテーション方式

- 15の**第一層ラベル (カテゴリ)** (例: 「生命兆候 Vital Signs」 「身体検査 Physical Examination」)
- 58の**第二層ラベル (詳細ラベル)** (例: 「血圧 Blood Pressure (BP)」 「脈拍 Pulse (HR)」)
- 階層的アノテーション:**
Step 1: 第一層ラベル (カテゴリ) を適用
Step 2: 第一層の範囲内で第二層ラベルを適用

第1層 15個	主訴 Chief Complaint (CC) 1 現在の病気の履歴 History of Present Illness (HPI) 2 既往歴 History 3 システムレビュー Review of Systems (ROS) 4 投薬/アレルギー Current Medications, Allergies 5 意識レベル Japan Coma Scale (JCS) 6 生命兆候 Vital signs 7 身体検査 Physical examination 8 実験検査 Laboratory examination 9 画像検査 Imaging examination 10 病理検査 Pathological examination 11 他の検査データ Other examination data 12
第2層 58個	病歴 Medical History 13 手術歴 Surgical history 14 家族歴 Family history 15 社会歴 Social History 16 体温 Temperature (BT) 17 脈拍 Pulse (HR) 18 血圧 Blood Pressure (BP) 19 呼吸 Respiratory Rate (RR) 20 経皮的動脈血酸素飽和度 Oxygen Saturation (SpO2) 21 体重 Body Weight (BW) 22 身長 Body Height (BT) 23 全身状態 Performance Status (PS) 24 血糖 Blood Sugar (BS) 25 全身所見 General Appearance 26 頭部 Head 27 目 Eyes 28 耳 Ears 29 鼻 Nose 30 口腔 Oral cavity 31 咽喉 Throat 32 頸部 Neck 33 心臓 Heart 34 肺 Lungs 35 腹部 Abdomen 36 四肢 Extremities 37 皮膚 Skin 38 神経系 Neurological Exam 39 リンパ節 Lymph Nodes 40 泌尿器 Genitourinary System 41 骨格系 Musculoskeletal System 42 血液検査 Blood Tests 43 尿検査 Urinalysis 44 免疫学的検査 Immunological Tests 45 微生物学的検査 Microbiological Tests 46 ホルモン検査 Hormonal Tests 47 血液ガス分析 Arterial Blood Gas (ABG) 48 腫瘍マーカー Tumor Markers 49 X線検査 X-ray Imaging 50 CTスキャン Computed Tomography (CT Scan) 51 MRI Magnetic Resonance Imaging 52 超音波検査 Ultrasound Imaging 53 核医学検査 Nuclear Medicine Imaging 54 造影検査 Contrast Imaging 55 内視鏡検査 Endoscopy 56 細胞診 Cytology 57 組織診 Histopathology 58 免疫組織化学 Immunohistochemistry (IHC) 59 分子病理 Molecular Pathology 60 特殊染色 Special Stains 61 病理解剖 Autopsy 62 呼吸機能 Respiratory Function 63 心電図 Electrocardiogram (ECG/EKG) 64 神経学的検査 Neurological Examinations 65 動脈血圧モニタリング Arterial Blood Pressure Monitoring 66 内服 Oral administration 67 点滴/注射 Intravenous administration 68 吸入 Inhalation administration 69 外用 Topical administration 70

Figure 2: Annotation Labels

3. モデル・手法

- LLaMA3.3: 70b
- Few-shot Prompting (5例) + Prompt Engineering + Retrieval-Augmented Generation (RAG)** (160例のゴールドスタンダードデータ)

Future

- 未登録ラベルが生成される問題 → **Promptの最適化が必要**
- 第二層ラベルの適用に課題あり → **Prompt Tuningが必要**
- 評価を行い
 - Precision (適合率)、Recall (再現率)、F1-score
 - LLMの出力と医師のアノテーションを比較し、精度を評価
 - 以下の2つの状況は正しいと判定
 - 完全に合わせる
 - 完全に合わせないが、80%以上のテキスト類似性を持つ場合

Results

✓ 良い結果

- Vital Signs (生命兆候) など、**定型的情報は比較的正確に分類される**
 - 例: 「SpO2 96%, BP 135/80, HR 80」は適切に「生命兆候」として分類
 - 第二層のラベル (「血圧」, 「脈拍」, 「SpO2」) も適切に付与
- 第一層 (カテゴリ分類) は概ね意図したラベルで識別される**
 - 「生命兆候 (Vital Signs)」, 「身体検査 (Physical Examination)」のカテゴリで特に高精度

```
{
  "text": "0) SpO2 94-96%, Bp146/96, P124, BT 37.8(",
  "labels": ["生命兆候 Vital signs"]
},
{
  "text": "SpO2 94-96%",
  "labels": ["経皮的動脈血酸素飽和度 Oxygen Saturation (SpO2)"]
},
{
  "text": "Bp146/96",
  "labels": ["血圧 Blood Pressure"]
},
{
  "text": "P124",
  "labels": ["心拍数 Heart Rate"]
},
{
  "text": "BT 37.8(",
  "labels": ["体温 Body Temperature"]
},
}
```

Figure 3: 良い結果の例 (モデルの一部出力)

△ 誤りの傾向

- 未登録ラベルの生成:** 事前に定義されていないラベルが生成された
 - 例: 「症状」や「赤血球数」など、定義されていないラベルが出現
- 第二層ラベルの適用が不完全**
 - 例: 「胸部写真: もともと右上肺野の透過性が低下しているが、その隙間の透過性が低下」が「画像診断」には分類されるが、第二層の「X線検査」のラベルが付与されなかった

```
{
  "text": "胸部写真: もともと右上肺野の透過性が低下しているが、その隙間の透過性が低下",
  "labels": ["画像診断 Imaging Studies"]
},
{
  "text": "A) 気管支肺炎の疑い",
  "labels": ["診断 Diagnosis"]
},
}
```

Figure 4: 第二層ラベルの適用が不完全の例

Discussion & Conclusion

本研究では、SOAPデータの階層的ラベリングをLLMを用いて自動化する試みを行った。結果として、第一層ラベルの適用は概ね良好であるが、第二層ラベルの適用にばらつきが見られた。特に以下の点が課題として浮かび上がった:

- 未登録ラベルの出力: **ほとんどは過剰推論のため、Prompt をさらに強化する必要あり**
- 第二層ラベルの適用が不完全: **Few-shotの例数が足りないのが原因と考えられるが、これ以上増やすと文字数が多すぎるため、Prompt Tuningが必要** これらの問題は、**モデルがラベルの適用ルールを完全には理解できていないこと、および入力テキストの構造的なばらつきに起因すると考えられる。**

[1] Linking clinical variables with health-related quality of life. A conceptual model of patient outcomes . I.B. Wilson et al. 1995
[2] LungDiag: Empowering Artificial Intelligence for Respiratory Disease Diagnosis Through Electronic Health Records. Hengrui Liang, Et al. 2023
[3] Large language models as data preprocessors . Zhang H . et al. 2023