

# ループリックに基づいたタグを付与した 日本語小論文データの構築と自動採点への効果

成岡 智也<sup>1</sup> 竹内 孔一<sup>2</sup>

<sup>1</sup> 岡山大学工学部 <sup>2</sup> 岡山大学大学院環境生命自然科学学域

pak13nrx@s.okayama-u.ac.jp

takeuc-k@okayama-u.ac.jp

## 概要

本研究では、小論文にループリックに基づいたアノテーションを行い、そのデータを用いて小論文自動採点を行う方法を提案する。ループリックの評価部分に対応した小論文の文書内構造の情報を、小論文にタグをつけるという形でアノテーションを行った。そうして得られるタグを小論文と同時にトークン化し、採点モデルにスコアを予測させ、小論文のみでの自動採点と比較することで、提案手法の有効性を明らかにした。

## 1 はじめに

大量の小論文を人間が平等に採点することは困難であり、またそれに伴う労力も大きい。近年、このような問題を解決するために、大規模言語モデルを用いた様々な小論文自動採点手法の研究が行われている。[1]しかし、小論文自動採点分野では、更なる精度の向上や、学習データを大量に必要にしてしまう点など様々な課題が見られる。

そこで、小論文を自動で採点するにあたって、追加情報を与えることを考える。小論文は、基本的に各問題に対して採点基準(ループリック)を基準に採点が行われる。小論文の内容に関して得点を決める際、その要素が小論文の中にどの程度含まれるかによって得点が決まる。先行研究では、小論文の構造を扱った研究がある[1][2][3][4]が、ループリックで評価される部分を日本語の小論文に対して付与したデータを作成し、小論文自動採点を適用した研究は我々の知る限り見当たらない。このような背景から、本論文ではループリックに基づいたアノテーションを行ったデータを作成するとともに、そのアノテーション結果を自動採点に利用し、小論文のみを利用した自動採点手法と比較し、有効性を明らかにする。

にする。

## 2 関連研究

大規模言語モデルによる小論文自動採点に小論文アノテーションを用いた例として、Kuらの研究が挙げられる[2]。この研究では、中国語小論文において、各小論文に対しての採点根拠を小論文外部にアノテーションし、それらを学習させることによって、モデルが未知の小論文に対してより良い採点根拠を出力させるを行っている。また、エンコーダに入力するトークンとして小論文以外のものを用いた例として、Chuらの研究が挙げられる[5]。この研究では、英語小論文において、GPT3.5 turboとLlama3.1-8B-Instructを用いて各小論文に対応する「根拠」を生成し、それぞれを共通のエンコーダに入力して自動採点を行うことによって、スコア予測において全体的なQWKの向上がみられている。本研究では、日本語小論文に対して内部にアノテーションを行い、その構造を追加で用いて自動採点を行っている。

## 3 ループリック評価部分のアノテーションデータ構築

本研究で使用している小論文データは、言語資源協会から公開されている日本語小論文データ[6]である。この小論文データは、2016年及び2018年に開講された講義の受講者の回答データから作成されている。今回利用した小論文データは「グローバル化の光と影」、「自然科学の構成と科学教育」、「東アジア経済の現状」、「批判的思考とエセ科学」の計4つのテーマであり、それぞれの講義テーマごとに3つの設問が用意されている。各設問ごとに300件前後の小論文がある。以降、各講義テーマは「グローバル化の光と影」をglobal,

表 1 文書構造分析タグの種類

課題名	タグ
global_q1	光, 影, 格差縮小, 格差拡大
global_q2	光, 影, 具体例, 企業名
global_q3	光, 影, 具体例, 見解
science_q1	実証性の説明, 再現性の説明, 客観性の説明
science_q2	自然相手, 持続役割, 根拠, 客観確保, 共通役割
science_q3	導入, 科学リテラシー, 展開, 今後
easia_q1	相互依存, 協力・協業の実態, 具体例
easia_q2	概略, 脱する方法
easia_q3	日本, 韓国, 中国, 協調と対立
criticize_q1	論理的・合理的思考の説明, 目標志向的思考の説明, 内省的・熟慮的思考の説明
criticize_q2	導入, 相関あり, 因果なし, 妥当でない理由, カラーテレビ, 平均余命, まとめ
criticize_q3	実例, 方法, 証拠, 説明, 要因

「自然科学の構成と科学教育」を science, 「東アジア経済の現状」を easia, 「批判的思考とエッセ科学」を criticize と表現する。加えて各設問を q1, q2, q3 と表現し、これらを組み合わせて「グローバル化の光と影」の設問 1 を global\_q1 という風に表現する。この各設問に対する各解答に、ループリックの内容をもとにしたタグを付けるアノテーションを行った。以降、このタグを「文書構造分析タグ」または単に「タグ」と呼ぶ。各設問に対応する文書構造分析タグの種類は、表 1 の通りとなっている。

global の設問および science の設問には、「光?」のようにタグに?がついているものも存在する。これは、同名のタグと比較して確信度が低い場合に用いられている。本研究では?がついているタグは一律に使用していない。

また、例えば設問 global\_q2 に関して、データ中のタグの種類は表に記載のある通り 4 種類であるが、うち「具体例」タグと「企業名」タグは理解力のスコアに関係がない。このような場合、一部タグは使用していない場合がある。

実際のアノテーションされた小論文は、図 1 のような形式になっている。

図 1 は、設問 global\_q1 の小論文の一部である。例えば、1 行目に着目すると、「グローバル化は世界全体の所得格差を縮小する」の部分に「光」「格差縮小」タグが、「国内及び各国間での所得格差を拡大している」の部分に「影」「格差拡大」タグがアノテーションされていることがわかる。

## 4 実験

本章では、3 章で述べたデータセットを用いて小論文の採点を行い、小論文採点における文書構造分析タグの有効性について確かめる。以下、各節で実験設定、実験結果、考察を述べる。また、本実験では、小論文採点モデルとして BERT ( Bidirectional Encoder Representations from Transformers ) [7] を使用する。

### 4.1 実験設定

実験では、BERT を用い、各小論文の 1 点から 5 点で評価されている「理解力」のスコアを 5 クラス分類によって予測する。本実験では、HuggingFace の BERT<sup>1)</sup>を利用する。評価指標には Accuracy と QWK ( Quadratic Weighted Kappa ) を用いる。それぞれ、以下の式で示す。

$$P(a, b) = \begin{cases} 1, & \text{if } a = b \\ 0, & \text{if } a \neq b \end{cases} \quad (1)$$

$$\text{Accuracy} = \frac{\sum_{l=1}^n P(A_l, B_l)}{n} \quad (2)$$

$$\text{QWK} = 1 - \frac{\sum_{i,j} O_{ij} W_{ij}}{\sum_{i,j} E_{ij} W_{ij}} \quad (3)$$

Accuracy は、式 (1), (2) によって求められる。n は各設問における小論文の総数であり、 $A_l$ ,  $B_l$  は l 番目の小論文に対してそれぞれ正解のスコア、予測したスコアである。Accuracy は 0 から 1 の値を取り、1 に近いほど一致率が高い。また、QWK は式 (3) によって求められる。i, j はそれぞれ正解のスコア、予測したスコアである。 $W_{ij}$  は正解のスコアと予測したスコアが離れているほど大きくなる値であり、観測された一致度  $O_{ij}$  と期待される一致度  $E_{ij}$  の双方に  $W_{ij}$  を重み付けする。QWK は、予測したスコアと正解のスコアが離れていればいるほど大きなペナルティを与える、小論文自動採点モデルで一般的に用いられる評価指標である。QWK は-1 から 1 の値を取り、1 に近いほど一致率が高い。

また、アノテーションをしてある小論文について、以下の図 2 ような方法を用いてトークナイズする。

これは、A というタグが 2 箇所、B というタグが 1 箇所についている小論文の例である。各タグに

1) <https://huggingface.co/tohoku-nlp/bert-base-japanese-v2>

格差縮小  
光

グローバリゼーションは世界全体の所得格差を縮小する一方で、国内及び各国間での所得格差を拡大している。これは、

光

グローバリゼーションに伴う自由主義貿易や経済支援によって、経済最下層の人々も含めた世界全体の利益上昇があった

格差拡大  
影

先進国での競争の活性化により発展途上国での利益を遥かに上回る利益を上げた結果、最恵国と最貧国あるいは経済力

図1 アノテーションされた小論文データ

**図2** BERTによる採点モデルと入力トークン例

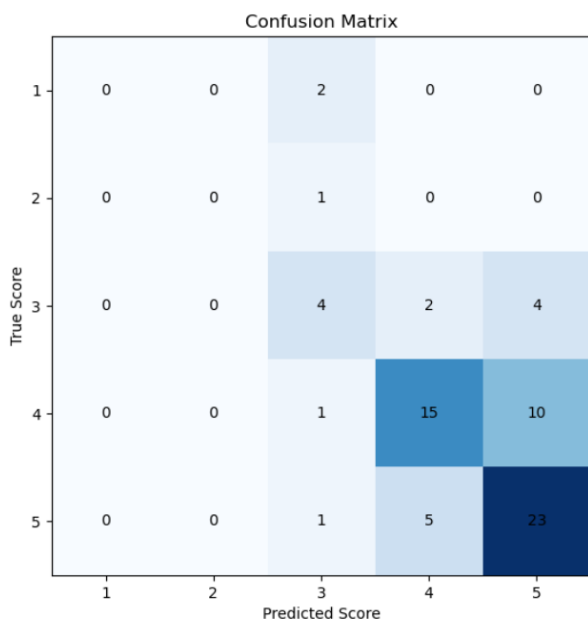


図3 science\_q1 タグ無し

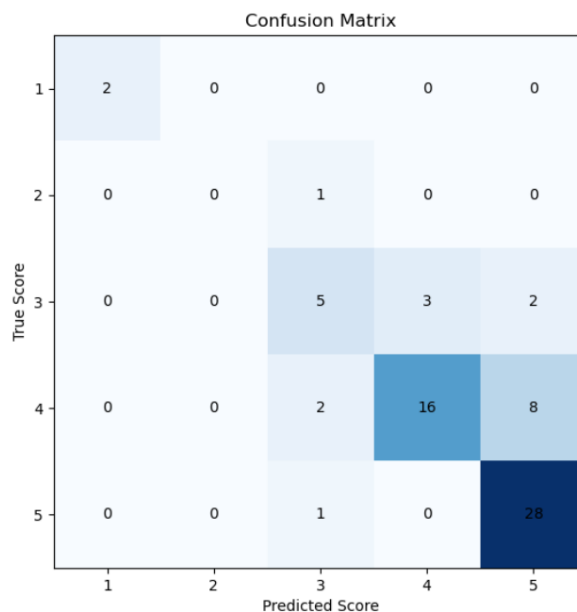


図4 science\_q1 タグ有り

## 5 まとめ

本研究では、BERT を用いた小論文自動採点において、小論文のアノテーションを行うことで得られる文書構造分析タグを用いる手法を提案し、その性能評価を行った。実験により、タグを用いることで、Accuracy と QWK の双方の向上が期待できることを明らかにした。しかし、一部の小論文課題に対してはスコア予測の性能の低下も見られたため、文書構造分析タグの利用方法についての課題があることも確認した。今後は、今回性能が低下した小論文課題に対しても有効に働くようなタグの利用方法の検討や、大規模言語モデルを用いたアノテーション及びタグ付けの自動化に取り組んでいく予定である。

## 6 データ

本研究で構築したデータは整理した後、公開する予定である<sup>2)</sup>。

2) [https://www.cl.cs.okayama-u.ac.jp/?page\\_id=2](https://www.cl.cs.okayama-u.ac.jp/?page_id=2) で案内する予定である。

の小論文と比較してその他の点数の母数が少なく、データセットに母数が少ないとタグがあっても予測が難しいといった側面は見られるため、こういった点は課題として挙げられる。

**性能が低下したものの考察** 設問 global\_q2, 設問 global\_q3, 設問 science\_q2, 設問 easia\_q2 では、タグを用いることで Accuracy と QWK のいずれか、もしくは双方において低くなってしまいう結果が得られた。しかし、設問 science\_q2 には大きな差があるわけではないことから、性能が落ちたのではなくタグがあまり採点に影響しなかったと考えるべきである。設問 global\_q3, 設問 easia\_q2 に関しては、QWK が上がってはいるものの Accuracy は低下しているため、性能が向上した、低下したと一概に評価することは難しい。設問 global\_q2 のみ Accuracy, QWK の双方において低下が見られる。これに関して、設問 global\_q2 が他の設問と比較して、Accuracy および QWK が高いことに起因しているのではないかと考えられる。このことから、タグを用いての小論文採点は小論文単体でスコアを評価することが簡単である設問には不適切である可能性が考えられる。これは、スコアを推定するための入力として小論文のみで十分であるにも関わらず、過剰な情報をタグによって与えてしまった結果、スコア予測に悪影響を与えてしまったということが想定される。

## 謝辞

本研究は JSPS 科研費 JP22K00530 の助成を受けたものです。

## 参考文献

- [1] 石井雄隆, 舟山弘晃, 松林優一郎, 乾健太郎. 国語記述式問題自動採点システムの開発と評価. 日本教育工学会研究報告集, 2024.
- [2] Nicoletta Calzolari, Min-Yen Kan, Veronique Hoste, Alessandro Lenci, Sakriani Sakti, and Nianwen Xue. Plaes: Prompt-generalized and level-aware learning framework for cross-prompt automated essay scoring. In **The International Conference on Computational Linguistics**, 2024.
- [3] 加藤嘉浩. 論理構造グラフを用いた自動採点モデル. 言語処理学会 第 30 回年次大会 発表論文集, 2024.
- [4] Yuning Ding, Marie Bexte, and Andrea Horbach. Score it all together: A multi-task learning study on automatic scoring of argumentative essays. In **Association for Computational Linguistics**, 2023.
- [5] SeongYeub Chu, JongWoo Kim, Bryan Wong, and MunYong Yi. Rationale behind essay scores: Enhancing s-llm’s multi-trait essay scoring with rationale generated by llms. In **arXiv:2410.14202**, 2024.
- [6] 竹内孔一, 大野雅幸, 泉仁宏太, 田口雅弘, 稲田佳彦, 飯塚誠也, 阿保達彦, 上田均. 研究利用可能な小論文データに基づく参照文書を利用した小論文採点手法の開発. 情報処理学会論文誌, 2021.
- [7] Jacob Devlin, Ming-Wei Chang, Kenton Lee, and Kristina Toutanova. Bert: Pre-training of deep bidirectional transformers for language understanding. In **Proceedings of the 2019 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies, Volume 1(Long and Short Papers)**, 2019.