

Slack ボットによる選好ラベル収集に基づく ニュース提示対話のスタイル採択予測

奥野 竜斗^{1,3} 吉野 幸一郎^{2,3} 飯尾 尊優^{1,3}

¹ 同志社大学 ² 東京科学大学

³ 理化学研究所ガーディアンロボットプロジェクト

ryuto.okuno@riken.jp tiio@mail.doshisha.ac.jp

koichiro.yoshino@riken.jp

概要

ニュースを話題とした人と対話システムの対話では、ユーザが求める支援の形と返答の方針が合わないと、受け取りにくさや違和感が生じ、満足度の低下や早期終了の一因となり得る。本研究は、ニュース提示対話における返答方針を三つのスタイルに整理し、文脈に応じてユーザがどのスタイルを選ぶかを予測することを目標とする。この目的のために、Slack 上でニュース候補の提示と三つの返答候補の提示を行い、ユーザの選択を正解ラベルとして記録するボットを開発してデータを収集した。また、収集したデータからユーザに提示するニュースの提示スタイルを予測するモデルを構築した。本提案システムでは日常的な対話環境で低負担に選好ラベル付きデータを収集し、スタイル採択予測の学習と評価に利用できる。

1 はじめに

ニュース提供は、新聞や放送による一方向の提示から、Web や SNS を通じた双方向のやり取りへと広がってきた。近年は音声やチャットによる対話形式も一般化し、ニュースについて質問しながら理解を深める利用形態が報告されている [1, 2, 3]。

一方で、ニュースは内容が多様であり、受け手の関心や感情も状況により変化する。同じニュースであっても、事実関係を整理したい場合もあれば、不安や怒りを受け止めてほしい場合もある。また、関連する背景を知りたい、別の視点を知りたいという場合もある。このような違いに対して、システムの返答方針を一つに固定すると、対話の流れに合わない返答が増え、対話が途切れやすくなる。

本研究は、ニュース提示対話における返答方針をスタイルとして明示化し、ユーザの発話や直前のやり取りに応じて、どのスタイルが選ばれやすいかを予測する枠組みを構築する。具体的には、Slack ボットにより選好ラベル付きの対話ログを収集し、そのデータを用いて複数の予測モデルを構築して比較評価する。あわせて、教師ありデータを継続的に収集するための実装とログ設計を示す。具体的には、1 回のターンごとに三つの返答候補を提示し、ユーザが最も適切と感じる候補を選ぶ仕組みを用いる。選好に基づく学習は、人の判断を学習信号として用いる方法として提案されている [4]。本研究は、この考え方をニュース対話のスタイル選択に適用し、収集の実装とログ設計を具体化する。

本稿の貢献は次の三点である。第一に、Slack 上でユーザの選択を正解ラベルとして記録できる収集基盤を設計し、実運用に近い操作で教師ありデータを取得する手順を示す。第二に、ニュース提示対話に必要な返答方針を三つのスタイルとして整理し、候補提示とラベル付けの単位を明確にする。第三に、収集ログを用いたスタイル採択予測の初期評価を行い、収集データが学習に利用可能であることを示す。

2 関連研究

2.1 対話データセットと返答方針

感情に配慮した対話の研究では、共感を含む応答を学習するためのデータセットが整備されている [5, 6]。また、外部知識に基づいて説明を行う対話として、参照知識を用いる枠組みが提案されている [7, 8]。しかし、ニュース提示では、説明の明確さ、

感情への配慮，理解を深める問い掛けを状況に応じて切り替え用いる必要がある。

ニュースを題材とした対話インタフェースの研究では，ニュースの提示方法がユーザの受け取り方に影響することが報告されている [1]。本研究は，提示されたニュースに対する対話の進め方に注目し，返答方針の選択を教師あり学習の対象として扱う。

2.2 評価信号としての行動ログ

対話システムの自動評価を人の満足度と一致させることは難しい [9]。対話システムの学習においては，ユーザが明示した選好や行動ログが自然に学習信号として蓄積される設計が重要である。情報検索分野では，クリックがフィードバックとして利用されてきた [10]。本研究では Slack ボットとの対話でユーザの選好が自然に収集される基盤をデザインする。

3 データ収集基盤

本章では，ニュースについて対話を行う機能と，その過程で教師ありデータを収集する機能を同時に満たす Slack ボットを述べる。ユーザはニュースを選び，返答候補から一つを選びながら対話を進める。この操作は対話の自然な流れとして実行でき，選択結果を正解ラベルとして蓄積できる。

実装には Slack を用いる。Slack は日常的に利用される環境であり，専用の実験用アプリを配布せずに運用できる。また，スレッド機能により複数の対話を同時に運用できる。

3.1 三つのスタイル定義

本研究は，ニュース提示対話における返答方針を三つに整理する。第一に論理である。論理は，ニュースの要点を整理し，因果や背景を明確にすることで理解を支援する。第二に共感である。共感とは，ユーザの感情を受け止め，安心感や受容を与えることを重視する。第三に探索である。探索は，追加質問や視点の提示により，ユーザの関心を深掘りし，理解を広げることを重視する。

この三分類は，ニュースを話題とした対話で求められやすい支援の形を，目的の違いに基づいて整理するために採用した。まず，知識に基づいて内容理解を助ける対話が研究されており [11]，ニュースの要点や背景を明確にする返答方針は重要になり得る。また，相手の感情に配慮した共感的応答を対象

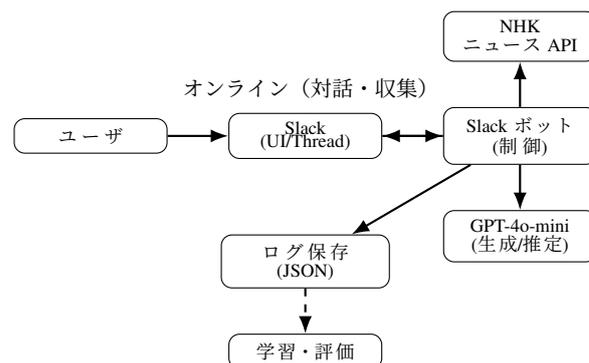


図1 システム構成と情報の流れ

とする研究もあり [12]，ニュースに対する受け手の感情を受け止める返答方針も必要になり得る。さらに，会話の中で追加質問を行い，情報の不足を補いながら理解を進める枠組みが提案されている [13]。以上を踏まえ，本研究では返答方針を論理，共感，探索の三つに整理し，候補提示とラベル付けの単位として用いる。

3.2 システム構成

図1 にシステム構成の概要を示す。ユーザは Slack 上でボットを呼び出し，ニュース選択と候補選択を行いながら対話を進める。ボットはニュース取得，候補生成，提示，ログ保存を担当する。外部 API から取得したニュース情報と，各ターンで提示した候補文，ユーザの選択結果を保存し，予測モデルの学習と評価に利用できる形で蓄積する。

処理は次の通りである。まず，NHK NEWS WEB の公開 API から最新ニュースを 100 件取得し，ランダムに 5 件を提示する。ユーザは一件を選択するか，ニュース選択を行わずに対話を開始する。ニュースを選ばない選択肢を用意する理由は，ニュースが関心に合わない場合や，ニュース以外の話題から入りたい場合に，無理な選択を避けるためである。

ボットは三つのスタイルに対応する返答候補を生成し，図2 に示すように三択として Slack 上に提示する。ユーザは最も適切と感じる候補を選ぶ。選択された候補は，次ターンでの対話履歴として採用される。この手順により，ユーザは自然な対話を行いながら，毎ターン一つの正解ラベルを付与する。

3.3 ログ設計

学習に用いるため，ログはターン単位で保存する。表1 は記録対象の主要項目である。ユーザの



図 2 ユーザが行う応答選択の例

表 1 ターンログ JSON の主要項目

項目	内容
timestamp	記録時刻
user_id	参加者識別子
thread_ts	セッション識別子
news_title, news_url	選択ニュースの題名と URL
candidates	三候補の本文と対応スタイル
selected_style	採択されたスタイル
selected_text	採択された候補文
user_input	ユーザの自由記述発話
skip_news	ニュース選択の有無

入力文、三つの候補文、選択結果を必ず含める。ニュースを選択した場合はニュースの題名と URL を保持する。この設計により、文脈と採択スタイルの対応を教師ありデータとして扱える。

ログの保存単位をターンにすることで、学習では一つのターンを一つの訓練データとして扱える。また、スレッド識別子により、セッション単位の分割や集計が容易になる。これは、評価時に同一セッションが学習と評価に同時に含まれる問題を避けるためにも有用である。

設計方針は次の五点である。第一に、ユーザ操作を最小化する。ユーザは候補から一つを選ぶだけでよい。第二に、提示順の影響を抑える。候補の並び順は毎回ランダムにし、位置による偏りを減らす [14]。第三に、情報を損なわない。採択された候補だけでなく、採択されなかった候補も保存し、学習と分析に利用できるようにする。第四に、再現可能な形で保存する。ニュースの URL などをログに保持し、後から参照できるようにする。第五に、個人情報の取り扱いを限定する。分析に必要な識別子は

表 2 収集データの概要

項目	値
参加者数	9
セッション数	18
ターン数 (正解ラベル数)	180
候補文総数	540
論理の採択率	24%
共感の採択率	33%
探索の採択率	43%

最小限にし、内容の再利用時に特定が起きにくい形式で保存する。

4 収集データ

表 2 に収集データの概要を示す。本稿では参加者 9 名から 18 セッション、180 ターンのログを収集した。各ターンでは三つの返答候補を提示しているため、候補文の総数は $180 \times 3 = 540$ である。このうち、各ターンでユーザが選択した候補が正解ラベルとなるため、正解データは 180 件である。一方で、選択されなかった候補も含めて保存されるため、候補生成の分析や、将来的な学習設計の拡張に利用できる。スタイル分布は論理が 24%、共感が 33%、探索が 43%であった。

4.1 収集の狙い

本基盤の狙いは、ユーザが自然に対話を続ける過程で、毎ターン必ずスタイル採択ラベルが得られるようにすることである。対話後のアンケートでは、回答者が場面を十分に思い出せない場合がある。また、評価に時間がかかり、長期的なデータ収集が難しくなる場合がある。本研究は、対話の最中に、その場面に合う候補を選んでもらうことで、文脈と選好を直接結び付ける。

この設計では、一つのユーザ発話に対して三つの候補が残る。そのため、採択された候補だけでなく、採択されなかった候補も合わせて保存できる。この情報は、後から候補の差分を分析する際に有用である。また、採択されにくい候補の特徴を整理することで、候補生成の改善にもつながる。

5 スタイル採択予測

5.1 予測タスク

収集ではユーザが望むニュース提供のスタイルをユーザが選択したが、将来的にはこの選択をシステ

ムが行うことが望ましい。そこで三値分類（チャンスレート 0.333）による予測を行う。今回限られた収集データを用いて予測可能性に関する基礎検討を行う。

5.2 比較条件

比較条件は次の四つである。

- **ゼロショット推定**：GPT-4o-mini に対して、三候補のうちどれが採択されるかを選ばせるための簡単な指示のみを与え、追加学習を行わずに推定させる。
- **ロジスティック回帰**：候補文と文脈から作成した特徴量を入力とし、L1 正則化付きロジスティック回帰で採択スタイルを予測する [15]。特徴量は、文脈文と各候補文の文埋め込み、それらの類似度、候補文の長さなどの指標を用いる。
- **ファインチューニング**：収集ログを教師信号として GPT-4o-mini をファインチューニングし、採択スタイルを直接予測する分類器として用いる。
- **perplexity 最小**：各候補の perplexity を計算し、値が最も小さい候補を採択候補として選ぶ。候補文を $y = (y_1, \dots, y_T)$ 、条件となる入力 x を連結し、応答生成を行う言語モデル (GPT-4o-mini) でトークン長で正則化した perplexity の計算を行う。

5.3 結果

表 3 に結果を示す。GPT-4o-mini のファインチューニングが 0.533 で最も高く、ロジスティック回帰もチャンスレートを上回った、ゼロショット推定と perplexity 最小はチャンスレートとの差がほとんどなかった。以上より、モデルが出力しやすい候補を選ぶだけでは、ユーザが適切と感じる候補を十分に捉えられないことが示唆される。これに対して、選好ラベルを用いたファインチューニングは精度が高く、ユーザの選択傾向を学習することで予測性能が改善する可能性がある。ただし、本稿のデータは 180 ターンのみであり、推定の安定性と一般性を評価するには追加収集が必要である。

表 3 スタイル採択予測の Top-1 Accuracy

方法	Accuracy
チャンスレート	0.333
ゼロショット推定: GPT-4o-mini	0.339
ロジスティック回帰	0.400
ファインチューニング: GPT-4o-mini	0.533
perplexity 最小	0.350

6 考察

6.1 収集基盤としての意義

本基盤の意義は、対話の自然な流れの中で、教師ありデータを継続的に収集できる点にある。対話システム評価では、自動指標のみでは実使用時の満足を説明しにくい [9]。本基盤は、ユーザがその場面で適切と感じた候補を直接ラベルとして蓄積するため、評価と学習を同一の操作で接続できる。また、採択されなかった候補も保存するため、候補間の差分分析や、誤りの傾向把握にも利用できる。

6.2 偏りと信頼性への配慮

提示順の影響は完全には除去できない。本稿では提示順のランダム化により偏りを抑えたが [14]、今後は選択確率の補正や、閲覧行動の測定も検討する必要がある [10]。

さらに、ニュース領域では事実と異なる内容が混入する危険がある [16]。本研究はニュースの URL をログに保持して再確認できる形を採用したが、将来的には参照情報に基づく返答生成の導入も重要になる [8]。本稿は収集基盤の提示を主目的とするため、事実確認の手順の詳細化は今後の課題とする。

7 結論

本稿では、ニュース提示対話におけるスタイル採択予測を目的として、Slack 上で選好ラベル付き対話ログを収集するボットを開発し、収集データに基づいて複数手法で予測を比較した。9 名から 18 セッション、180 ターンのログを収集し、各ターンで三候補を提示することで候補文総数 540 件を得た。初期評価では、GPT-4o-mini によるゼロショット推定、ロジスティック回帰、GPT-4o-mini のファインチューニング、perplexity 最小を比較し、ファインチューニングが最も高い精度を示した。今後は、データ規模と多様性の拡大、提示順の影響低減、およびセッション単位での厳密な評価を進める。

8 謝辞

本研究は JST RISTEX JPMJRS24L3 の支援を受けた。

参考文献

- [1] Harshita Sahijwani, Jason Ingyu Choi, and Eugene Agichtein. Would you like to hear the news? investigating voice-based suggestions for conversational news recommendation. In *Proceedings of the 2020 Conference on Human Information Interaction and Retrieval*, pp. 437–441. ACM, 2020.
- [2] Koichiro Yoshino and Tatsuya Kawahara. News navigation system based on proactive dialogue strategy. In *Natural language dialog systems and intelligent assistants*, pp. 15–25. Springer, 2015.
- [3] Hiroaki Takatsu, Katsuya Yokoyama, Yoichi Matsuyama, Hiroshi Honda, Shinya Fujie, and Tetsunori Kobayashi. Recognition of intentions of users’ short responses for conversational news delivery system. In *INTERSPEECH*, pp. 1193–1197, 2019.
- [4] Paul F. Christiano, Jan Leike, Tom Brown, Miljan Martić, Shane Legg, and Dario Amodei. Deep reinforcement learning from human preferences. In *Advances in Neural Information Processing Systems (NeurIPS)*, 2017.
- [5] Hannah Rashkin, Eric Michael Smith, Margaret Li, and Y-Lan Boureau. Towards empathetic open-domain conversation models: A new benchmark and dataset. In *Proceedings of the 57th Annual Meeting of the Association for Computational Linguistics*, pp. 5370–5381, Florence, Italy, 2019. Association for Computational Linguistics.
- [6] Yukun Ma, Khanh Linh Nguyen, Frank Z. Xing, and Erik Cambria. A survey on empathetic dialogue systems. *Information Fusion*, Vol. 64, pp. 50–70, 2020.
- [7] Emily Dinan, Stephen Roller, Kurt Shuster, Angela Fan, Michael Auli, and Jason Weston. Wizard of Wikipedia: Knowledge-powered conversational agents. In *International Conference on Learning Representations (ICLR)*, 2019.
- [8] Patrick Lewis, Ethan Perez, Aleksandra Piktus, Fabio Petroni, Vladimir Karpukhin, Naman Goyal, Heinrich Küttler, Mike Lewis, Wen-tau Yih, Tim Rocktäschel, Sebastian Riedel, and Douwe Kiela. Retrieval-augmented generation for knowledge-intensive NLP tasks. In *Advances in Neural Information Processing Systems (NeurIPS)*, 2020.
- [9] Chia-Wei Liu, Ryan Lowe, Iulian Serban, Mike Noseworthy, Laurent Charlin, and Joelle Pineau. How NOT to evaluate your dialogue system: An empirical study of unsupervised evaluation metrics for dialogue response generation. In *Proceedings of the 2016 Conference on Empirical Methods in Natural Language Processing*, pp. 2122–2132, Austin, Texas, 2016. Association for Computational Linguistics.
- [10] Thorsten Joachims, Laura A. Granka, Bing Pan, Helene Hembrooke, and Geri Gay. Accurately interpreting click-through data as implicit feedback. In *Proceedings of the 28th Annual International ACM SIGIR Conference on Research and Development in Information Retrieval (SIGIR)*, pp. 154–161, Salvador, Brazil, 2005. ACM.
- [11] Emily Dinan, Stephen Roller, Kurt Shuster, Angela Fan, Michael Auli, and Jason Weston. Wizard of wikipedia: Knowledge-powered conversational agents. *arXiv preprint arXiv:1811.01241*, 2018.
- [12] Hannah Rashkin, Eric Michael Smith, Margaret Li, and Y-Lan Boureau. Towards empathetic open-domain conversation models: A new benchmark and dataset. In *Proceedings of the 57th annual meeting of the association for computational linguistics*, pp. 5370–5381, 2019.
- [13] Mohammad Aliannejadi, Hamed Zamani, Fabio Crestani, and W Bruce Croft. Asking clarifying questions in open-domain information-seeking conversations. In *Proceedings of the 42nd international acm sigir conference on research and development in information retrieval*, pp. 475–484, 2019.
- [14] Nick Craswell, Onno Zoeter, Michael Taylor, and Bill Ramsey. An experimental comparison of click position-bias models. In *Proceedings of the International Conference on Web Search and Web Data Mining (WSDM)*, pp. 87–94. ACM, 2008.
- [15] Nils Reimers and Iryna Gurevych. Sentence-bert: Sentence embeddings using siamese bert-networks. In *Proceedings of the 2019 Conference on Empirical Methods in Natural Language Processing and the 9th International Joint Conference on Natural Language Processing (EMNLP-IJCNLP)*, pp. 3982–3992, Hong Kong, China, 2019. Association for Computational Linguistics.
- [16] Ziwei Ji, Nayeon Lee, Rita Frieske, Tiezheng Yu, Dan Su, Yan Xu, Etsuko Ishii, Yejin Bang, Andrea Madotto, and Pascale Fung. Survey of hallucination in natural language generation. *ACM Computing Surveys*, Vol. 55, No. 12, pp. 248:1–248:38, 2023.