

対話システムに任意のスタイルを付与する生成ベースのスタイル制御手法の提案

Proposal of a generated base style control method
to apply arbitrary styles to dialogue systems京江 遼大[†] 打矢 隆弘[†] 内匠 逸[†]
Ryota Kyoe Takahiro Uchiya Ichi Takumi

1. はじめに

近年、対話システムが普及し、単に人間の対話相手になるだけでなく、人間のように多種多様なスタイルを持って対話そのものを目的とする高度な対話システムの需要が高まっている [1]. 特に、ニューラルネットワークを用いた生成ベースの対話システムが数多く提案されている。しかし、このような対話システムでは、様々な話者の発話を基にした大規模な対話コーパスを学習データに用いるため、学習したモデルが生成する発話において話者のスタイル一貫性に欠けることが多いという欠点が存在する [2]. 例えば、「私」や「僕」等の一人称が混在した文章を生成する可能性がある。

本研究では、文面に現れる一人称・三人称や語尾・方言などが話者ごとに異なると仮定し、話者の特徴的な話し方をスタイルと定義し、一貫性に欠けた発話を減らし、対話システムの開発者が意図したスタイルを簡単かつ自由に付与することを目的とする。スタイルを学習する対話モデルに与え、その学習結果を基に、ユーザとの対話が可能な対話システムの実現を目指す。モデルの応答文にスタイルを反映させるための単純な手法として、スタイルを自然言語で直接モデルの入力文に与える手法 [3] が考えられる。しかし、この手法では与えるスタイル情報が増えれば増えるほど入力文が長くなり、モデルの入力上限を超える等の問題があるため、適した方法であるとは言えない。また、スタイル情報を直接入力文に与えるため、スタイル情報に登場する単語に強く条件付けられ、前後の文脈とは関係なしにそのまま反映された応答文が生成され、対話として不自然なものになる可能性も考えられる。そこで、入力文とは別にスタイル情報を埋め込む手法を提案する。具体的には、漫画のキャラクターのセリフをスタイル情報として扱い、漫画の対話をデータセットとして用いる。また、対話応答の自然さを向上させるために、対話システムの応答前の数ターンの対話履歴を文脈情報として埋め込むことも行う。各埋め込みを学習によって最適化することで、対話性とスタイル性を両立させた応答文の生成を目指す。

本研究では、人手評価による被験者の主観評価を行う。対話システムにスタイルが反映できているか、応答が自然であるか、文法的に正しいか、これらの項目を 5 段階で評価する。

2. 関連研究

2.1 対話システムのスタイル制御

対話システムが人間とさらに自然な対話を行うためには、

[†]名古屋工業大学大学院工学研究科

Nagoya Institute of Technology, Graduate School of Engineering

一貫した話者であること、知識を持つこと、感情を持ち対話相手に共感すること、の 3 要素が必要である [4]. 本研究では、話者のスタイルを扱うことで、一貫した話者であることに注力する。

ニューラル対話モデルにおいて初めて応答文に現れる話者のスタイルの制御を目指した Li らの研究 [2] では、ある話者の大量の発話から話者の情報を分散表現として獲得し、デコーダでの応答文生成の際に利用する **Speaker Model** を提案した。このモデルでは話者 ID を付与した対話データセットを使用して、単語の分散表現と同様の形で話者の分散表現を獲得する。そして、話者の分散表現と入力文を連結し、応答文生成デコーダへの入力とする。

また、Zhang らの研究 [3] では、話者の情報を明示的に記述したプロフィール文を含んだ対話データセットである **Persona Chat** を用いて、話者の情報を応答文へ反映させるモデルを提案した。このデータセットは、5 文程度のプロフィール文で構成されたセットを作成し、これに則った対話をクラウドワーカーにさせる事で収集したものである。話者の情報を記述した文章とその話者の対話データを用いて、プロフィール文と入力文を与えられると、話者の情報に基づいた応答文を生成可能なモデルである。

2.2 Prompt Tuning

BERT[5] や GPT[6] 等の事前学習済みモデルである大規模言語モデルの登場により、自然言語処理での深層学習を用いる場合は、大規模なデータで事前学習したモデルを **Fine-Tuning** し、目的の言語タスクに適応させる手法が現在では主流となった。しかし、言語モデルはさらに大規模化が進み、**Fine-Tuning** の学習コストの増大が問題になり、大規模言語モデルの持つ膨大な知識を活用したパラメータを更新せずに目的タスクへ適応させる新たな手法 **Prompt-Tuning** が注目を浴びている。多くの場合はモデルに入力する文章を **Prompt** として扱い、モデル全体のパラメータを更新するのではなく、与える文章に変化を与えることで目的のタスク達成を行う手法である。

Prompt-Tuning の一例として、Brow らの研究 [7] では、タスクの説明といくつかのサンプルを **Prompt** として作成し、目的のタスクを実行する **Zero/Few-Shot Learning** を提案した。これらの **Prompt** は手作業で数例を作成することで効果が得られたが、**Fine-Tuning** よりも精度は悪いことが知られている。作成例を図 1 に示す。

また、**Prompt** を最適化する手法の中には、2.1 節で述べたような方法な **Prefix-Tuning**[8] がある。入力の前頭に **Prefix-Token** と呼ばれるトークン列を追加し、事前学習済みモデルの各パラメータは更新せず、このトークン列のパラメータのみを最適化することで目的のタスクを解く。**Prefix-Tuning** の概要を図 2 に示す。

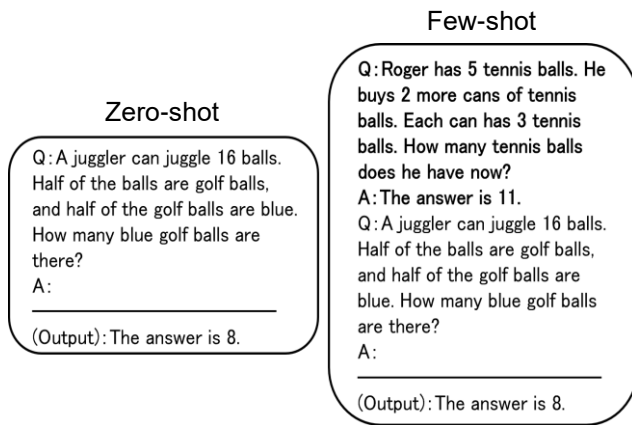


図 1 Zero-shot と Few-shot のサンプル

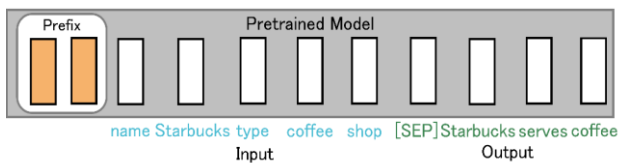


図 2 Prefix-Tuning の概要図

3. 提案手法

3.1 提案モデル

大規模言語モデルの知識を扱うために、Transformer ベースの事前学習済みモデルに、スタイル情報を埋め込む Style Embedding 層、対話履歴を埋め込む Dialog Embedding 層を新たに追加したモデルの構築を提案する。この提案モデルのアーキテクチャを図 3 に示す。

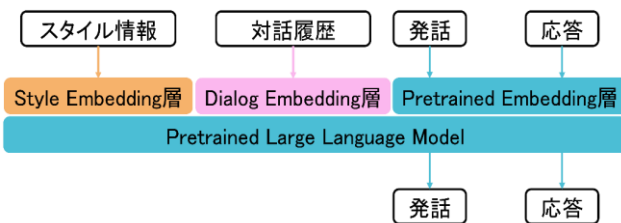


図 3 提案モデルの全体図

3.2 使用データセット

本研究では、目的のために話者のスタイルを多分に含んだデータセットを使用する必要がある。そこで、漫画のキャラクターが多く収集されている Manga109 データセット [9] を用いる。Manga109 データセットは 1970 年代から 2010 年代の漫画 109 冊を基にしたものである。漫画と対応したアノテーションがされており、テキストではセリフ、モノローグや状況説明をアノテーション対象として収録している。データセットの構築では、セリフをスタイル情報として扱うために、キャラクターごとにランダムに 20 発話を選択した。これは発話数が少ないキャラクターを除外するためである。また、対話履歴として数ターンの会話を用意するために、最大 5 ターンの会話を選択した。大きなターン数を扱うことは文脈を捉えることで重要であるが、今回は prompt として入力するため膨大な長さにならないかつ、漫画は展開が激しく動くものであるため、多くのターンを扱うよりは数ターンに抑えた。

3.3 学習方法

スタイル情報である漫画のキャラクターセリフを Style Embedding 層、キャラクター同士の対話履歴を Dialog Embedding 層、該当対話ターンの発話と応答のペアを事前学習済み言語モデルの Embedding 層によって、各々埋め込みを行う。これらの埋め込みベクトルを結合し、事前学習済み言語モデルに入力する。大規模言語モデルをベースにしているため、Fine-Tuning の学習コストが高いため、学習時には、新たに追加した Style Embedding 層と Dialog Embedding 層のパラメータのみ更新を行う。

4. 実験

3.3 節の手法による、スタイルを応答文に反映できる対話システムの構築を行う。モデルの作成には Hugging Face の Transformers を使用し、学習環境には Google Colaboratory を利用した。学習に使用した GPU は NVIDIA A100 であり、GPU メモリは 40GB である。また、提案モデルとの比較のために、事前学習済みモデルをデータセットの発話ペアで Fine-Tuning するベースラインモデルの作成も行う。3.2 節のデータセットは日本語の漫画データセットのため、日本語事前学習済みモデルの rinna GPT[10]を利用した。

4.1 評価

各モデルを組み込んだ対話システムの生成文が発話に対する応答として、自然で話者の一貫したスタイルが反映されているかどうかを人手評価する。評価方法は Zhang らの手法[3]を参考に、応答が文法的に正しいか、対話として自然か、話者のスタイルが反映されているか、の 3 項目を 5 段階で評価し、各評価者の平均を最終的な評価値として扱う。

5. まとめ

話者の特徴をスタイルとして扱う対話システムのスタイル制御手法の提案を行った。今後は、人手評価により提案手法の有効性を検証する。

参考文献

- [1] 総務省, “高度対話エージェント技術の研究開発の推進”, 令和元年版情報通信白書第 2 部, pp.387-388, 2019.
- [2] Li Jiwei et al, “A persona-based neural conversation model”, arXiv preprint arXiv:1603.06155, 2016.
- [3] Zhang Saizheng et al, “Personalizing Dialogue Agents: I have a dog, do you have pets too?”, arXiv:1801.07243, 2018.
- [4] Stephen Roller et al, “Recipes for building an open domain chatbot”, arXiv:2004.13637, 2020.
- [5] Jacob Devlin et al, “BERT: Pre-training of Deep Bidirectional Transformers for Language Understanding”, arXiv:1810.04805, 2018.
- [6] Radford A et al, “Improving Language Understanding by Generative Pre-Training”, Technical Report OpenAI, 2018.
- [7] Tom B. Brown et al, “Language Models are Few-Shot Learners”, arXiv:2005.14165, 2020.
- [8] Xiang Lisa Li and Percy Liang, “Prefix-Tuning: Optimizing Continuous Prompts for Generation”, arXiv:2101.00190, 2021.
- [9] “Manga109 dataset,” <http://www.manga109.org>.
- [10] 趙 天雨, 沢田 慶, “日本語自然言語処理における事前学習モデルの公開”, 人工知能学会研究会資料 言語・音声理解と対話処理研究会 93 巻, pp169-pp170, 2021.