

対話システムにおけるペルソナ更新の導入とペルソナの数による影響分析

吉田快¹ 品川政太郎^{1,2} 須藤 克仁^{1,2} 中村 哲^{1,2}

¹ 奈良先端科学技術大学院大学

² 理化学研究所革新知能統合研究センター

{yoshida.kai.yf1, sei.shinagawa, sudoh, s-nakamura}@is.naist.jp

概要

既存のペルソナ対話システムのペルソナは初めに条件付けられたもので固定されており、自動的に更新することができない。そのため、応答に含まれる新たなペルソナを考慮できないという課題がある。本研究ではシステムの応答履歴に応じて自動的にペルソナを更新するという新しい問題設定を考え、これを実現するために、ペルソナ追加機構を持つペルソナ対話システムを提案し、その際に起こりうる課題の調査を行った。

1 はじめに

ペルソナと呼ばれる複数の記述文によって表現されるプロフィール情報をシステムに持たせることで、そのペルソナに沿った一貫性のある応答をする対話システムとしてペルソナ対話システム [1] がある。ペルソナ対話システムに対する既存のアプローチでは、いくつかの文を明示的なシステムのプロファイルとして組み込むことが試みられている [2]。しかし、既存のアプローチではペルソナが初めに条件付けられたもので固定されており、自動的に更新することができない。例えば、図 1 のように、システムに条件付けられていないペルソナに関しての質問(図 1: ユーザー入力 1)に対しては、“I have two dogs” (システム応答 1) のような新たなペルソナを含んだ応答を生成する可能性がある。そのため、システムが一貫した対話を行うにはこのような新たなペルソナも考慮する必要がある。そこで、本研究では、ペルソナ対話システムが応答履歴に応じて自身のペルソナを自動で更新するという新たな問題設定を考え、これを実現するために、ペルソナ追加機構を持つペルソナ対話システムを提案する。ペルソナ追加機構はシステムの応答履歴から Mazare らの研

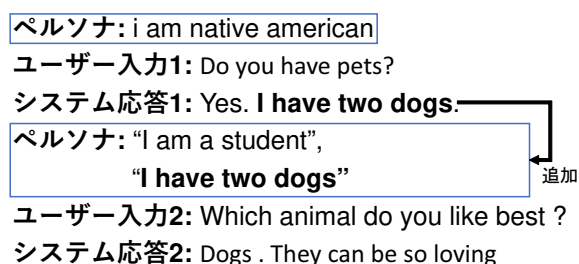


図 1 ペルソナ追加機構を持つ対話システム (提案手法) を用いた会話例

究 [2] と同様に条件を満たす記述をペルソナとして抽出し、抽出されたペルソナをフィルタリングして既存のペルソナに追加する機構である。

提案手法では、図 1 のシステム応答 1 のように新しいペルソナを含む応答をした場合に、条件を満たす記述を新しいペルソナとして追加し、次の応答生成に用いることで、ユーザー入力 2 に “Dogs. They can be so loving.” と応答することが期待できる。提案手法では既存のペルソナに新しいペルソナを随時追加するため、ペルソナの数の変化が応答生成に影響を及ぼす可能性がある。そこで本研究では、この影響について調査を行い、ペルソナを更新する問題設定における課題について明らかにする。

ペルソナ追加機構のフィルタリングには、入力された 2 つのペルソナ文が含意・中立・矛盾のいずれの関係にあるか分類する分類器 [3] を用いる。しかし、文献 [3] では含意のみしか扱っていない。そこで本研究では含意・中立・矛盾を含む評価用データセットである Persona NLI データセットを新たに構築し、分類器の精度評価を行うことで、この分類器の精度が本研究に用いる上で十分か確かめる。

2 関連研究

明示的なペルソナを用いた応答生成に、含意関係推論を組み合わせることで、限られた規模の対話

データで対話モデルを構築できる BERT-over-BERT (BoB) [4] がある。BoB では、ペルソナ対話データで学習された応答生成用のデコーダと、非対話型推論データで学習された、生成した応答をペルソナ文により一貫した表現に再調整するデコーダの2つのデコーダを用いることで、少量の対話データでのペルソナを含んだ応答生成を可能にしている。しかし、BoB は事前に与えられたペルソナでシステムのペルソナが固定されているため、長い対話を想定した場合、実際に生成された応答と異なった応答をすることが考えられる。

一方で、Transformer を用いて特定のユーザーの過去の対話履歴から暗黙的なペルソナを自動的に学習する手法 [5] がある。この手法はユーザーの応答履歴を逐次更新するため、実質的にペルソナの更新を応答ごとに行っているが、ペルソナが暗黙的であることから、ペルソナを最大限活用できていないことが考えられる。

本研究では、応答履歴から逐次的に明示的なペルソナを更新することで、長い対話における応答の矛盾を削減し、よりペルソナを活用できるシステムが構築できることを期待する。

3 手法

3.1 応答生成モデル

本研究では、ペルソナ追加機構の有無による応答生成への影響を検証するために、ベースとなる応答生成モデルとして、BERT-over-BERT (BoB) [4] を用いる。今回は著者らが公開しているコード¹⁾を用いた。BoB は明示的なペルソナを用いた応答生成に、含意関係推論を組み合わせることで、限られた対話データで学習を可能にするモデルである。BoB の構造を図 2 に示す。BoB は図に示すようにエンコーダ E と応答デコーダ D_1 , D_2 の3つの BERT ベースのサブモジュールで構成されている。 D_1 はペルソナ P とクエリ Q から応答ベクトル生成を行うデコーダである。また、 D_2 は NLI コーパスによって、与えられたペルソナ P に対して、より一貫した応答 R_2 を生成するように訓練されたデコーダであり、 D_1 の出力である応答ベクトルとペルソナ P から応答を再調整する。これにより、 D_2 は D_1 よりもさらにペルソナに対して一貫した応答を生成することができる。生成された応答候補群からサンプリングする方

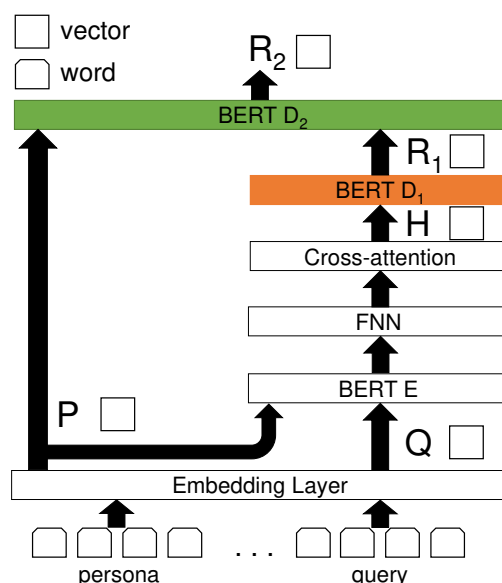


図 2 BoB の構造

法には、ペルソナの更新による影響を確認するため先行研究と同様の greedy search を用いた。

3.2 ペルソナ追加機構

システムの応答履歴からペルソナを更新するために、ペルソナ追加機構を導入する。ペルソナ追加機構は、システムの応答履歴 H から新しいペルソナを抽出し、既存のペルソナに追加するために2つの段階を踏む。まず初めに、(a) 応答履歴 H から Mazareらの研究 [2] と同様に以下の3つの条件を満たす記述をペルソナとして抽出する。

1. 4 から 20 の単語 (句読点を含む) で構成される
2. “I” か “my” という単語を含む
3. 名詞、代名詞、形容詞のうち少なくともどれか1つを含む

そして、(b) 抽出されたペルソナを NLI モデルを用いた重複、矛盾フィルタの2種類のフィルタによって更新の有無を決定する。NLI モデルは、入力された2つの文の関係を含意・中立・矛盾の3つのラベルに分類する分類器である。本研究における NLI モデルは、事前学習済みの BERT を、ペルソナ対話データに基づいて構築された Dialogue NLI コーパス [3] によって fine-tuning することで作成できる。本研究では、公開されている学習済みモデル [3] を用いた。このモデルのペルソナの組み合わせにおける中立ラベルの予測精度は 94.48% とされている [3] が、含意や矛盾ラベルにおける予測精度は示されていないため、独自に Persona NLI データセットの構

1) <https://github.com/songhaoyu/BoB>

築を行い、精度評価を行う。

4 実験設定

実験では、まず検証に用いるモデル (BoB) を作成し、ペルソナ追加機構に用いる NLI モデルの精度を調査する。これらを踏まえて、ペルソナの数による応答生成への影響について議論する。

4.1 モデル構築

BoB は学習済みモデルを公開しておらず、新たに学習する必要がある。本研究では、応答生成モデルの構築のために、ペルソナが付与された対話データセットである Persona-Chat [1] と、2つの文章間の含意関係を推測するタスクの、非対話型推論データセットである MNLI コーパス [6] を用いた (データセットの詳細は付録 A に記載した)。

4.2 NLI モデルの精度調査

ペルソナの含意関係推論における既存の NLI モデルの精度調査を行うため、新たに評価用のデータセットを構築し、このデータセットに対する NLI モデルの分類結果を評価する。本節では、この評価用のデータセットである Persona NLI データセットの構築方法について述べる。

Persona NLI データセット

Persona NLI データセットはペルソナのペアに対して、MNLI コーパスと同様に含意 (e)、中立 (n)、矛盾 (c) をラベル付けしたデータセットである。データの総数は 401 であり、内訳は含意が 103、中立が 196、矛盾が 102 となっている。データは次のようにペルソナ文のペアを作り、ラベル付けを行った。

含意 Persona-Chat からペルソナをランダムに選択し、含意になるようなペルソナを人手で生成した。例えば、Persona-Chat に含まれている “i love watching superheroes shows” に対して、それと含意関係になるような “i like superheroes” を人手で生成しペアを作った。

中立 Persona-Chat に定義されているペルソナの中から、ペルソナを中立になるように選択し組み合わせた。

矛盾 Persona-Chat からランダムに選択したペルソナに対し、2通りの方法で矛盾したペルソナのペアを作成した。1つ目は、“i work in a veterinary office” を “i work in supermarkets” のように矛盾し

た文に書き換える方法であり、2つ目は否定を用いて “i love beef” を “i don't like beef” のように書き換える手法である。52 件を 1つ目の方法で、50 件を 2つ目の方法で作成した。

4.3 ペルソナの数による応答生成への影響の調査

ペルソナの数を変化させ、応答生成への影響を調べる方法について述べる。まず、定義されているペルソナの数である 5 個より少ない場合は不必要な数を削除した。6 個から 20 個の場合は重複や矛盾が無いように Persona-Chat に定義されているペルソナからランダムに追加し、20 個より多い場合は重複しないようにランダムに追加した。その後、数を変化させたペルソナと Persona-Chat の入力文を用いて応答生成を行った。

評価指標

ペルソナ対話システムにおいては、(a) 多様性があり、(b) ペルソナに一貫した応答を行うことが望ましい。そのため、ここでは応答の多様性と、応答とペルソナの一貫性の 2つの側面で生成された応答を評価する。応答の多様性は、distinct 1/2 (**Dist.1/2**) [7] を用いて評価を行う。Dist は式 (1) のように、生成された応答に含まれるユニークな n-gram の数を単語数でスケールした数によって、どの程度の語彙を含んだ応答なのか評価するものである。

$$\text{distinct-n} = \frac{|\text{unique}(n\text{-gram})|}{|\text{uni-gram}|} \quad (1)$$

応答の一貫性の評価には、Consistency Score (**C.Score**) [8] を用いる。これはモデルを用いてペルソナと生成した応答の一貫性を予測するもので、式 (2) のように生成された応答と応答生成に使用したペルソナを NLI 分類器によって含意、中立、矛盾に分類し、そのスコアの総和で評価を行う。

$$NLI(r, q_i) = \begin{cases} -1 & r \text{ が } p_i \text{ と矛盾} \\ 0 & r \text{ が } p_i \text{ と中立} \\ 1 & r \text{ が } p_i \text{ と含意} \end{cases}$$
$$C.Score(r) = \sum_{i=1}^l NLI(r, p_i) \quad (2)$$

5 実験結果

5.1 NLI モデルの精度

予測結果の混同行列を表 1 に示す。分類の結果、

表 1 分類結果

		予測		
		含意	中立	矛盾
正解	含意	60	32	11
	中立	1	190	5
	矛盾	46	10	46

全クラスの予測精度は約 74%であった。ラベル別に見ると、中立の正解率は約 97%と高い結果となったが、含意は 58%、矛盾は 45%という結果となった。矛盾を含意と誤って予測した原因として、語彙がかなり似ている文同士が含意と分類されやすいことが確認できた。例えば、“i like apple”と“i hate apple”は矛盾であるが、含意と予測されてしまった。

5.2 ペルソナの数による影響

はじめに、学習時と同じペルソナ 5 個 (再現) の場合とペルソナの更新 (提案手法) を導入した場合の評価を表 2 に示す。結果、Dist と C.Score の両方で提案手法が既存手法の再現を上回る結果であった。

表 2 再現したモデルと提案手法の評価

Filter	Dist.1	Dist.2	Dist.AVG	C.Score
再現結果 [4]	2.573	18.81	10.69	6.796
提案手法	3.161	21.34	12.25	16.50

次にペルソナの数ごとのスコアを図 3、図 4 に示す。結果として、C.Score はペルソナの数が増えるほどスコアが下がっていくことが確認できた。また、Dist についてはペルソナの数が多いほどスコアが高く、増えるにつれてスコアが下がることが確認できた。生成した応答の内容を確認してみると、ペルソナの数が増えるほど、入力に関わらず、ペルソナのみに依存した応答をすることが確認できた。

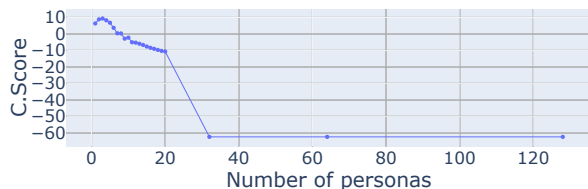


図 3 ペルソナの数による応答の一貫性への影響

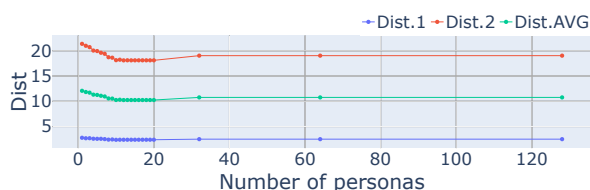


図 4 ペルソナの数による応答の多様性への影響

6 考察

表 2 の再現したモデルと提案手法では、システムの応答履歴からペルソナを増やすことで、Dist と C.Score の両方を増加させることが確認できた。一方でペルソナの数を変化させた応答生成においては、ペルソナが増えるにつれて Dist と C.Score は下がることが確認できている。原因として、ペルソナの数による実験では Persona-Chat に定義されている語彙を多く含むペルソナを追加しているため、モデルが機能し、複数のペルソナに沿った応答ができていたため Dist が下がったことが考えられる。一方で、提案手法の実験では応答からペルソナを抽出し追加しているため、Persona-Chat のペルソナに含まれていない語彙を含んだペルソナが追加される可能性がある。そのためモデルが機能せず Dist が上がっていると考えられる。実際に学習データに定義されているペルソナのサブワードの種類が 3,940 であるのに対し、応答生成に使用したペルソナの内 1,567 が学習データに含まれていない語であった。また、ペルソナの数が増えるまで Dist が低下し、それ以降はまた増加するのに関しても、同様にモデルが機能する限界と考えることができる。次に、ペルソナが 64 個以上の場合にスコアが同じになってしまった現象については、ペルソナ 1 つに含まれる単語数がおおよそ 6, 7 個であることから、64 個より大きいペルソナを設定した場合に、エンコーダの上限に達してしまっただと考えられる。C.Score については、ペルソナが増えることで一貫した応答ができないことが確認できた。しかし 20 個以上の場合は、ペルソナをランダムで選択しているため、その影響もあることが考えられる。

7 まとめ

本研究では、システムの応答履歴でペルソナを更新することで応答生成に与える影響を分析した。ペルソナの数による応答生成への影響の実験から、学習データに含まれていない語彙を持ったペルソナに対してはモデルが機能しないことや、ペルソナの数が増えたと応答が固定化されることが示唆された。そのため、今後の課題として、未知語に対応でき、かつ応答生成に必要なペルソナのみを選択できるモデル構築が必要だと考えられる。

参考文献

- [1] Zhang et al. Personalizing dialogue agents: I have a dog, do you have pets too? **arXiv preprint arXiv:1801.07243**, 2018.
- [2] Mazaré et al. Training millions of personalized dialogue agents. **arXiv preprint arXiv:1809.01984**, 2018.
- [3] Welleck et al. Dialogue natural language inference. **arXiv preprint arXiv:1811.00671**, 2018.
- [4] Song et al. BoB: BERT over BERT for training persona-based dialogue models from limited personalized data. In **Proc. ACL-IJCNLP, 2021**.
- [5] Zhengyi Ma et al. One chatbot per person: Creating personalized chatbots based on implicit user profiles. In **Proc. SIGIR, 2021**. ACM, 2021.
- [6] Williams et al. A broad-coverage challenge corpus for sentence understanding through inference. **arXiv preprint arXiv:1704.05426**, 2017.
- [7] Li et al. A diversity-promoting objective function for neural conversation models. **arXiv preprint arXiv:1510.03055**, 2015.
- [8] Lin et al. Personalizing dialogue agents via meta-learning. **arXiv preprint arXiv:1905.10033**, 2019.
- [9] Dinan et al. The second conversational intelligence challenge (convai2). **arXiv preprint arXiv:1902.00098**, 2019.

A 付録

Persona-Chat

Conv AI2 Persona-Chat (Dinan et al., 2019) [9] は、大規模なタスク指向ではない大規模なデータセットで、1,155 人の個性的なキャラクターが含まれており、それぞれが少なくとも 5 つのプロフィール文で構成されている。このデータセットは、Zhang らの作成した Persona-Chat (Zhang et al., 2018) [1] を応答候補が複数になるよう拡張したものである。Persona-Chat は Amazon Mechanical Turk で収集されており、ペアの話者はそれぞれ与えられたプロフィールを条件に対話を行ったデータが収録されている。

本研究ではこのコーパスを 121,880 の学習用に、9,558 を検証用に、7,504 をテスト用に分割して使用した。

MNLI コーパス

Multi-Genre Natural Language Inference (MNLI) コーパス (Williams et al., 2018) は文章理解のための機械学習モデルの開発・評価に使用するために設計されたデータセットで、およそ 43 万件の自然言語推論 (文の含意関係の認識) のために利用可能なコーパスである。

データの内訳は、含意が 130,615、中立が 130,590、矛盾が 130,590 となっている。