

# 対話における重要語抽出を用いた文脈と解釈の逐次相互依存解決

吉野哲平<sup>1\*</sup> 松森匠哉<sup>1</sup> 福地庸介<sup>1</sup> 滝本佑介<sup>1</sup> 阿部佑樹<sup>1</sup> 新行内浩輔<sup>1</sup> 今井倫太<sup>1</sup>

<sup>1</sup> 慶應義塾大学 理工学部

**Abstract:** 自然言語処理において、発話文を理解することは解決すべき課題の一つである。連続的な発話文を扱うためのモデルとして Encoder-Decoder モデルや、SCAIN がある。本研究は、SCAIN のフレームワークに重要語抽出を組み込んだモデルを提案する。発話文の中の着目すべき単語を逐次推定することにより、既存の SCAIN より高い精度での文脈推定を実現した。また、重要語抽出においても既存手法を上回る性能を示した。

## 1 序論

全文が揃っている文章と比較して、発話文は対話者が交代で発話するため、その都度処理する逐次的処理が必要である。また、指示語や多義語といった曖昧性の高い単語を多く含んでいるため、処理がより困難である。したがって、発話文を適切に処理するためには単語がもつ語義の曖昧性を解消しながら文脈を逐次的に推定する必要がある。

語義の曖昧性解消や、文脈の逐次推定には、既存研究がある。語義の曖昧性解消の既存手法である、周辺単語を加味した単語埋め込みを行うモデル [1, 2] では、単語埋め込み表現をその単語に一意に定められた点とせず、文脈を加味した上での点として扱うことにより、出現文に合わせた単語解釈を行うことができる。また、語義の曖昧性を扱う別の手法として、単語の意味の広さを表現可能な単語埋め込みモデル [3] が存在する。通常の単語埋め込み [4, 5] では単語を点として表現するところを、[3] では点でなく正規分布で表すことにより、その単語の意味の広さや単語間の包含関係を表すことができる。また、文脈の逐次推定と語義の曖昧性解消が可能なモデルとして、SCAIN [6] がある。SCAIN は、発話単語の意味決定には文脈を考慮する必要があり、また文脈を決定するためには発話単語の解釈が必要であるという、文脈と単語解釈の相互依存を逐次的に解決するアルゴリズムである。SCAIN では、パーティクルフィルタにより文脈を同時に複数仮定し、単語をそれぞれの文脈上で解釈することができる。これにより、発話文が含む曖昧性を逐次的に解消しながら単語解釈を進めることができる。

しかし、これらの既存モデルは全て発話理解には不十分なものである。まず曖昧性を解消した単語埋め込みの手法は、複数の文脈を保持することができないために、発話への応用が困難である。既存の単語埋め込

み手法の曖昧性解消は、推定される単一の文脈をもとに行われている。しかし、発話においては文脈が一意に定められないまま対話が進むため、複数の文脈を保持しながら単語をそれぞれの文脈において解釈する必要があるが、既存の単語埋め込み手法はいずれも文脈の曖昧性を考慮した逐次解釈を行うことができない。一方で SCAIN は、複数の文脈のもとで発話単語を逐次的に解釈することが可能であるため、発話文における単語の曖昧性を解決できる。しかし SCAIN は、文脈に寄与する単語が発話文の中に占める割合が少ない場合に正しく文脈を推定できない点で問題があり、実験室的でない発話を解釈することは困難である。

本論文では、SCAIN に重要語抽出を組み込むことで、重要な単語を推定しながら発話文を解釈する SCAIN with keyword extraction を提案する。SCAIN のフレームワークに従うことで文脈と解釈の相互依存を逐次的に解決可能であり、複数文脈を同時に保持することができるという点で、曖昧性を排除した単語埋め込みよりも発話解釈に適している。また、重要語抽出を取り入れることで既存の SCAIN がもつ制約を解決した。本提案における重要語抽出は、SCAIN において文脈がパーティクルフィルタで表現されていることを利用し、それぞれの推定文脈に関連が強い語を重要語として扱うものである。これは、文章においてタイトルや概要を文脈保証の拠り所として重要語抽出を行うこと [7] と相似である。

本論文の構成は次のとおりである。まず提案の章で SCAIN with keyword extraction についてそのアルゴリズムを説明し、どのように問題が解決可能であるかを説明する。次に実験と評価の章で、本手法が獲得する文脈表現を既存の SCAIN モデルのものと比較し、また獲得した重要語についても既存手法と比較し、それぞれ評価する。

\*連絡先：慶應義塾大学

E-mail: yoshino@ailab.ics.keio.ac.jp

## 2 提案

本論文では、SCAIN の処理の内部に重要語抽出を取り入れた、SCAIN with keyword extraction を提案する。本手法は、分散表現空間上において単語と文脈の位置に関して、それらを相互依存的に解決するアルゴリズムである。

以下では、SCAIN が行う処理について説明する。本手法はシステムとユーザの発話対一つに対し、3つの処理を行う。Step1 では、エージェントが発話を行い、文脈を移動させる。Step2 では、ユーザが発話を行い、文脈を移動させる。ただし、ユーザの発話文に含まれる単語は解釈が一意に定まらず曖昧性をもつので、これについて各推定文脈に基づき解釈を行う。また、重要語抽出もそれぞれの推定文脈に基づいて行う。Step3 では、文脈候補の妥当性を評価し、文脈候補を伴うパーティクルを削減する。重要語や他の発話単語を参考に、各文脈の尤もらしさを算出し取捨選択する。

## 3 実験と考察

### 3.1 実験

SCAIN with keyword extraction により得られる発話文の解釈の妥当性を確かめるため、以下でその評価実験を行う。本提案手法の発話文解釈に関する出力は、文脈と重要度付き解釈地図の二つである。よって、以下では文脈と重要語抽出の二点を既存手法と比較しながら評価する。文脈に関しては、既存研究にて行われた実験と同様に、定量的評価も行う。

#### 3.1.1 文脈推定の評価

文脈の精度を評価するため、既存の SCAIN と本提案手法 SCAIN with keyword extraction に同一の文を入力し、文脈の移動を観察する。入力文として、“If you will play latest games, powerful computer will be better for delivering the power you need” を入力した。結果は図 1,2 のとおりである。

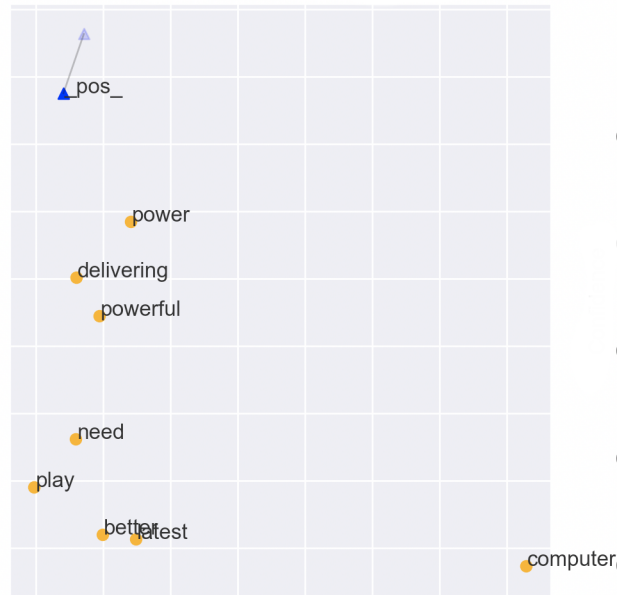


図 1: SCAIN による文脈推定

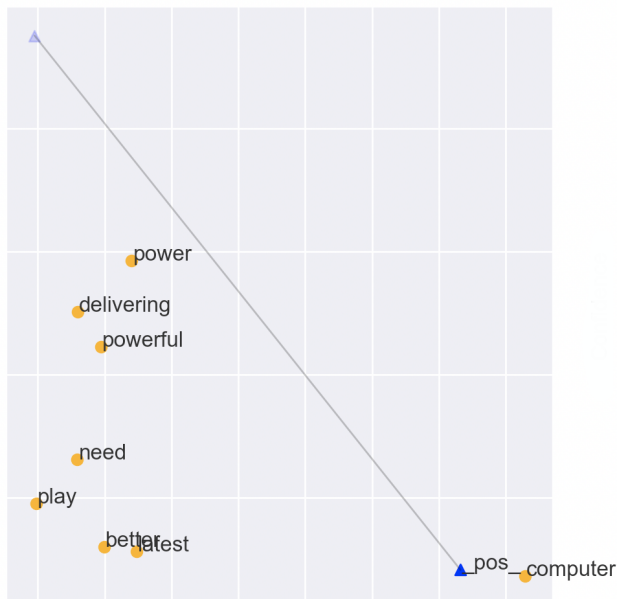


図 2: SCAIN with keyword extraction による文脈推定

ただし、pos が文脈の位置を表している。また、示されている平面は SCAIN が単語を管理する 100 次元空間を PCA を用いて 2 次元平面へと次元削減した物である。既存の SCAIN は文脈が図左上から微量移動したのに対して、SCAIN with keyword extraction では文脈が computer の方向へ移動した。

### 3.1.2 重要語抽出の評価

次に, SCAIN with keyword extraction の重要語抽出について他のアルゴリズムと比較する. 比較対象として, tf-idf[8] と RAKE[9] を用いる. 既存の SCAIN は設計上, 重要語抽出を想定していないため, ここでは比較対象から除く. また tf-idf では単一の文に対して重要語抽出を行うことができず, 複数の文と比較する必要があるため, 比較対象として nltk の NPS Chat Corpus を利用した. 本データセットはチャットルームにおける 10,000 以上の投稿をまとめたものである. 3.1.1 と同じく, “If you will play latest games, powerful computer will be better for delivering the power you need” を入力し, 重要語抽出の挙動を比較した. 結果は表 1 のとおりである.

表 1: 重要語抽出の比較 (値は相対的重要度)

| keyword SCAIN     | tf-idf            | RAKE                      |
|-------------------|-------------------|---------------------------|
| computer(0.165)   | delivering(0.158) | play latest games (0.529) |
| better(0.126)     | powerful(0.158)   | powerful computer (0.235) |
| latest(0.118)     | latest(0.137)     | power(0.059)              |
| games(0.115)      | games(0.137)      | need(0.059)               |
| need(0.108)       | power(0.137)      | delivering(0.059)         |
| power(0.107)      | computer(0.088)   | better(0.059)             |
| play(0.095)       | play(0.068)       |                           |
| powerful(0.088)   | better(0.060)     |                           |
| delivering(0.079) | need (0.057)      |                           |

提案手法では, computer や latest という単語が重要語であると推測しており, また delivering は重要語ではないと推測している. tf-idf では, 提案手法とは反対に, delivering が最重要単語であると推測している. また RAKE では, play latest games と powerful computer をイディオムとして認識しており, それらが重要語であるとしている.

## 3.2 考察

### 3.2.1 文脈推定の評価

図 1, 図 2 に関して文脈の動きを比較すると, 既存手法の結果である図 1 では文脈に寄与すべきでない語に引かれて, 文脈の移動が意味を成していないのに対し, 提案手法の結果である図 2 では重要語 computer を感知し文脈が computer へ移動していることがわかる. 発話文は書き言葉に比べ文脈への貢献が少ない単語を多く含むため, 文脈に寄与しない単語を判別して文脈の推定を行うことは重要である.

### 3.2.2 重要語抽出の評価

表 1 に関して重要語抽出を比較する. まず, 提案手法と tf-idf の結果を比較する. 特に, delivering の重要度が提案手法では最低値である一方, tf-idf では最高値である点で大きな差異が見られた. 入力文の内容から, 明らかに重要な語はコンピュータやゲームに関する語であり, delivering は重要語ではないため, tf-idf による重要語抽出は適切に働いていないことがわかる. tf-idf は単語頻度と逆文書頻度の積であるが, 発話は繰り返し表現に関してはほぼ全て代名詞を使用してしまうため, 単語頻度を重要度の尺度にすることはできない. つまり, 発話文に対する tf-idf は逆文書頻度による希少語フィルタとしてのみ機能する. このことは tf-idf の重要語抽出の結果を見ても明らかである.

次に, 提案手法と RAKE の結果を比較する. どちらも computer や game 周辺概念を重要語として観測しており, これらの結果は正しい. イディオムや係受けの情報を必要とするのであれば RAKE が適しており, 出力を単語として得たい場合は本提案手法が適している. NLP モデルにおいて, 単語の重要度を入力の補助情報を使用することは一般的であり, 特に Attention などの形でしばしばモデルに組み込まれる. 一方で, ランダムな N 単語セットの重要度を入力として扱うことは容易ではないため, 提案手法の出力の方が他の NLP モデルとの親和性は高い.

## 4 結論

SCAIN with keyword extraction は, SCAIN に重要語抽出を組み込むことで, 発話文の中の着目すべき単語を逐次推定することにより, 既存の SCAIN より高い精度での文脈推定を実現した. また, 重要語抽出においても既存手法を上回る性能を示した.

## 参考文献

- [1] Matthew E Peters, Mark Neumann, Mohit Iyyer, Matt Gardner, Christopher Clark, Kenton Lee, and Luke Zettlemoyer. Deep contextualized word representations. *arXiv preprint arXiv:1802.05365*, 2018.
- [2] Jacob Devlin, Ming-Wei Chang, Kenton Lee, and Kristina Toutanova. Bert: Pre-training of deep bidirectional transformers for language understanding. *arXiv preprint arXiv:1810.04805*, 2018.

- [3] Luke Vilnis and Andrew McCallum. Word representations via gaussian embedding. *arXiv preprint arXiv:1412.6623*, 2014.
- [4] Tomas Mikolov, Ilya Sutskever, Kai Chen, Greg S Corrado, and Jeff Dean. Distributed representations of words and phrases and their compositionality. In *Advances in neural information processing systems*, pages 3111–3119, 2013.
- [5] Jeffrey Pennington, Richard Socher, and Christopher D Manning. Glove: Global vectors for word representation. In *Proceedings of the 2014 conference on empirical methods in natural language processing (EMNLP)*, pages 1532–1543, 2014.
- [6] Yusuke Takimoto. Slam-inspired simultaneous contextualization and interpreting for conversation sentences. Keio University, 2019.
- [7] Rekha Bhowmik. Keyword extraction from abstracts and titles. In *IEEE SoutheastCon 2008*, pages 610–617. IEEE, 2008.
- [8] Daniel Jurafsky and James H. Martin. *Speech and Language Processing: An introduction to natural language processing, computational linguistics, and speech recognition*. 2000.
- [9] Stuart Rose, Dave Engel, Nick Cramer, and Wendy Cowley. Automatic keyword extraction from individual documents. *Text mining: applications and theory*, 1:1–20, 2010.