

RNN 系列変換モデルを用いた高階論理式からの文生成

馬目華奈 (指導教員: 戸次大介)

1 はじめに

近年の構文解析と意味解析の技術の発展によって、文の意味を論理式で表して高度な推論を行うシステムの構築が可能となった。このようなシステムは、含意関係認識 [1, 2] や文間類似度計算 [3] のタスクで高精度を達成しており、今後、さらなる自然言語処理タスクへの応用が期待されている。

文からその論理式への変換が高精度に行われる一方で、論理式を自然言語文に戻す方法については自明ではない。しかし、論理式から自然言語文に変換することができれば、推論システムの改善や、様々な自然言語処理タスクへの応用が期待できる。

そこで本研究では、機械翻訳等の系列変換において高い精度を示しているニューラルネットによる系列変換モデル (Sequence-to-Sequence model) [4] を用いて高階論理式から文を生成する手法を提案する。

論理式の埋め込みについては複数の方法を提案・比較した。含意関係認識用データセットを用いて提案手法の評価を行った結果、論理式を先頭の記号から埋め込んだ場合と比較して、論理式を木構造として埋め込むことで精度向上がみられた。

2 背景

2.1 CCG に基づく論理式による文の意味表現

文を高階論理式に変換し、高階論理に基づく自動推論を行うシステムとして、ccg2lambda [2] がある。ccg2lambda では、まず入力文に対して組合せ範疇文法 (Combinatory Categorical Grammar, CCG) [5, 6] に基づく統語解析を行う。CCG は語彙化文法の一つで、統語構造に並行して意味表示の合成を行う文法体系として知られている。各語には統語範疇が割り当てられる。CCG では語と語の統語的・意味的な関係を、関数適用や関数合成などの組合せ規則により計算していく。同時に文における各語の貢献の仕方を定めていく。

次に、CCG の導出木をラムダ計算に基づいて論理式へと変換する。ラムダ項によって表現されている各語の意味表示から、組合せ規則が指定する計算に沿って、最終的な文の意味表示である高階論理式を得ることができる。高階述語論理による証明には、高階論理・型理論に基づく定理証明支援系である Coq が用いられている。

2.2 系列変換モデル

系列変換モデル [4] とは入出力がシーケンスとなる機構で、意味や構文などには注目せず、入力と出力の対応を学習して覚える、ニューラルネットのモデルである。系列変換モデルは入力列を隠れ状態ベクトルに変換するエンコーダと、隠れ状態ベクトルから出力を行うデコーダからなる。

エンコーダでは、入力の系列を埋め込みベクトルに変換した後、LSTM 等の再帰型ニューラルネットワークによって隠れ状態ベクトルに変換する。デコーダでは、エンコーダで出力された隠れ状態ベクトルを初期

値とし、隠れ状態と自身のこれまでの出力結果を基に次のトークンを生成する。

3 提案手法

3.1 モデルとデータセット

本研究では意味表現の一つである、抽象的意味表現 (Abstract Meaning Representation, AMR) からの文生成においても高精度を実現している系列変換モデルを用いて、ccg2lambda が生成する高階論理式を入力とし、対応する自然言語文を予測することを目的とする。エンコーダ、デコーダには LSTM を用いた系列変換モデルを用いた。

学習データは、ccg2lambda を用いて文から高階論理式を作成する。実験用テキストには、含意関係認識タスクの評価用データセットである SNLI [7] を用いて、論理式と自然言語文のペアからなる教師データを作成した。CCG パーザの解析失敗を避けるため、SNLI データセットの文例のうち、一文に含まれる単語数が 60 単語以内の文例 20,000 件を対象とした。このうち、構文解析、論理式への変換に成功したユニークなデータ 12925 件を使用する。データのうち、9,140 件を教師データ (うち 2,285 件を validation データ)、1,500 件をテストデータとする。

3.2 論理式の埋め込み

“Bob ate pizza.” という例は、論理式に変換すると以下ようになる。

$$\text{exists } x.((x = \text{Bob}) \ \& \ \text{exists } z1.(\text{Pizza}(y) \ \& \ \text{exists } e.(\text{eat}(e) \ \& \ (\text{Subj}(e) = x) \ \& \ (\text{Obj}(e) = y))))$$

論理式を系列変換モデルの入力として扱うためには論理式をベクトル化する必要がある。論理式をベクトル化する上で、本研究では記号単位、およびトークン単位で数値化する手法を検討する。記号単位の場合は、論理式を先頭の記号から one-hot ベクトルに変換して埋め込みを行う。また、トークン単位で埋め込む場合は、論理式の構成要素を埋め込む順序も問題となる。本研究では先頭から順に埋め込む手法と、論理式を木構造やグラフとみなし、先行順 (pre-order) による深さ優先探索でリスト化する手法を採用した。論理式を木構造とグラフ構造で表した例を図 1 に示す。論理式の変数のノードを統一することでグラフ化を行う [8]。4 種類の埋め込み手法の例を下記に示す。(ここで $_$ はスペース文字を表す)。

1. 記号単位で先頭から埋め込む手法
例: [e,x,i,s,t,s,_,x,_,.,_,(,(x,_,=,...]
2. トークン単位で先頭から埋め込む手法
例: [exists,x,(,(x,=,Bob,)&,exists,...]
3. 木構造で埋め込む手法
例: [exists,x,&,,=,Bob,x,exists,y,&,...]
4. グラフで埋め込む手法
例: [exists,x,x,x,&,Bob,=,&,Pizza,...]

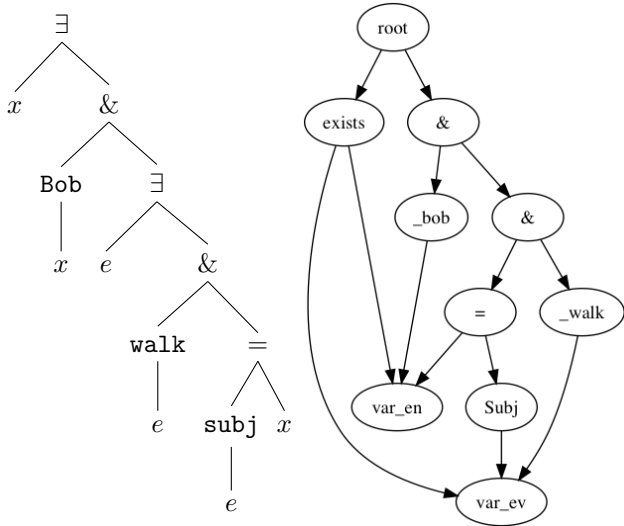


図 1: 木構造 (左) とグラフ構造 (右)

4 実験

4.1 実験設定

各埋め込み手法における入出力の語彙、文の最大の長さを表 1 に示す。

表 1: 各埋め込み手法における語彙数・系列長

	記号	トークン	木構造	グラフ
入力語彙数	70	5,118	5,107	4,991
出力語彙数	78	7,214	7,214	7,214
入力列最長	2,097	699	451	259
出力列最長	270	55	53	53

4.2 評価方法

評価は BLEU [9] によって行う。BLEU とは、文と文の表層的な類似度 (n -gram 一致率) を表す指標である。本研究では、以下の式に基づいて評価を行う。

$$score = BP \exp \left(\sum_{i=1}^N \frac{1}{N} \log P_n \right)$$

$$BP = \begin{cases} 1 & (c \geq r) \\ \exp \left(1 - \frac{r}{c} \right) & (c < r) \end{cases}$$

$$P_n = \frac{\sum_{i=0} \text{出力文 } i \text{ 中と解答文 } i \text{ 中で一致した } n\text{-gram 数}}{\sum_{i=0} \text{出力文 } i \text{ 中の全 } n\text{-gram 数}}$$

4.3 実験結果・考察

実験の BLEU スコアを表 2 に示す。グラフにして埋め込む手法が文の表層的な一致率が高かった。例として、“Two surgeons are having lunch.” という文が文がデコードされた結果を表 3 に示す。“surgeons” に対して各モデルが異なる代替となる語を出力していることがわかる。このような入力系列には存在しない単語を出力してしまうケースがあり、このような問題に対しては今後系列変換モデルのコピー機構 [10] を取り

入れることを検討する。記号単位以外のものについては、“having” も一致させることができていた。

表 2: 評価結果

指標	記号	トークン	木構造	グラフ
BLEU	34.9	39.7	41.8	44.7

表 3: 生成された文の例

文	Two surgeons are having lunch.
記号単位	Two children are playing basketball.
トークン単位	Two entertainers are having fun.
木構造	Two teams are having a brawl.
グラフ	Two brothers are having a picnic.

5 おわりに

本研究では、系列変換モデルを用いて高階論理式から文を生成する手法を提案した。含意関係認識用データセットを用いて提案手法の評価を行った結果、論理式をシーケンス化して先頭から埋め込んだ場合と比較して、論理式の順番を考慮して埋め込むことで精度向上がみられた。今後の課題として、他の意味表現からの文生成との比較や他のデータセットによる評価を行う。また、アテンション付き系列変換モデルやコピー機構を用いるなど、モデルの改良に取り組む。

参考文献

- [1] Lasha Abzianidze. A Tableau Prover for Natural Logic and Language. In *Proc. of EMNLP*, 2015.
- [2] Pascual Martínez-Gómez, Koji Mineshima, Yusuke Miyao, and Daisuke Bekki. ccg2lambda: A Compositional Semantics System. In *Proc. of ACL System Demonstrations*, 2016.
- [3] Hitomi Yanaka, Koji Mineshima, Pascual Martínez-Gómez, and Daisuke Bekki. Determining Semantic Textual Similarity using Natural Deduction Proofs. In *Proc. of EMNLP*, 2017.
- [4] Ilya Sutskever, Oriol Vinyals, and Quoc V Le. Sequence to Sequence Learning with Neural Networks. In *Proc. of NIPS*, 2014.
- [5] Mark Steedman. Surface Structure and Interpretation. In *The MIT Press*, 1996.
- [6] Daisuke Bekki. *A Formal Theory of Japanese Grammar: The Conjugation System, Syntactic Structures, and Semantic Composition*. Kuroshio, 2010. (In Japanese).
- [7] Samuel R. Bowman, Gabor Angeli, Christopher Potts, and Christopher D. Manning. A large annotated corpus for learning natural language inference. In *Proc. of EMNLP*, 2015.
- [8] Mingzhe Wang, Yihe Tang, Jian Wang, and Jia Deng. Premise selection for theorem proving by deep graph embedding. In *Proc. of NIPS*, 2017.
- [9] Kishore Papineni, Salim Roukos, Todd Ward, and Wei-Jing Zhu. BLEU: a Method for Automatic Evaluation of Machine Translation. In *Proc. of ACL*, 2002.
- [10] Jiatao Gu, Zhengdong Lu, Hang Li, and Victor O.K. Li. Incorporating Copying Mechanism in Sequence-to-Sequence Learning. In *Proc. of ACL*, 2016.